



City Research Online

City, University of London Institutional Repository

Citation: Jahromizadeh, Soroush (2013). Joint rate control and scheduling for providing bounded delay with high efficiency in multihop wireless networks. (Unpublished Doctoral thesis, City University London)

This is the unspecified version of the paper.

This version of the publication may differ from the final published version.

Permanent repository link: <https://openaccess.city.ac.uk/id/eprint/2454/>

Link to published version:

Copyright: City Research Online aims to make research outputs of City, University of London available to a wider audience. Copyright and Moral Rights remain with the author(s) and/or copyright holders. URLs from City Research Online may be freely distributed and linked to.

Reuse: Copies of full items can be used for personal research or study, educational, or not-for-profit purposes without prior permission or charge. Provided that the authors, title and full bibliographic details are credited, a hyperlink and/or URL is given for the original metadata page and the content is not changed in any way.

Joint Rate Control and Scheduling for Providing Bounded Delay with High Efficiency in Multihop Wireless Networks

By: Soroush Jahromizadeh
Supervisor: Dr Veselin Rakočević

Systems and Control Research Centre
Department of Electrical and Electronic Engineering

A Thesis Submitted to the School of Engineering and
Mathematical Sciences, City University London,
in Fulfilment of the Requirements for the Degree of
Doctor of Philosophy in Electrical Engineering

June 12, 2013



Abstract

This thesis considers the problem of supporting traffic with elastic bandwidth requirements and hard end-to-end delay constraints in multi-hop wireless networks, with focus on source transmission rates and link data rates as the key resource allocation decisions. Specifically, the research objective is to develop a source rate control and scheduling strategy that guarantees bounded average end-to-end queueing delay and maximises the overall utility of all incoming traffic, using network utility maximisation framework. The network utility maximisation based approaches to support delay-sensitive traffic have been predominantly based on either reducing link utilisation, or approximation of links as M/D/1 queues. Both approaches lead to unpredictable transient behaviour of packet delays, and inefficient link utilisation under optimal resource allocation. On the contrary, in this thesis an approach is proposed where instead of hard delay constraints based on inaccurate M/D/1 delay estimates, traffic end-to-end delay requirements are guaranteed by proper forms of concave and increasing utility functions of their transmission rates. Specifically, an alternative formulation is presented where the delay constraint is omitted and sources' utility functions are multiplied by a weight factor. The alternative optimisation problem is solved by a distributed scheduling algorithm incorporating a duality-based rate control algorithm at its inner layer, where optimal link prices correlate with their average queueing delays. The proposed approach is then realised by a scheduling algorithm that runs jointly with an integral controller whereby each source regulates the queueing delay on its paths at the desired level, using its utility weight coefficient as the control variable. Since the proposed algorithms are based on solving the alternative concave optimisation problem, they are simple, distributed and lead to maximal link utilisation. Hence, they avoid the limitations of the previous approaches. The proposed algorithms are shown, using both theoretical analysis and simulation, to achieve asymptotic regulation of end-to-end delay given the step size of the proposed integral controller is within a specified range.

Contents

1	Introduction	10
1.1	Multi-hop Wireless Networks	10
1.2	Multi-hop Wireless Networks Design	12
1.2.1	Physical Layer	12
1.2.2	Access Layer	14
1.2.3	Network Layer	16
1.2.4	Transport Layer	20
1.2.5	Application Layer	20
1.2.6	Motivations for Cross-Layer Design	21
1.3	Research Objectives	25
1.3.1	Background on QoS-Oriented Load Control	25
1.3.2	Problem Description and Assumptions	26
1.3.3	Limitations of Current Solutions	27
1.3.4	Research Objectives	27
1.4	Summary of Contributions	28
1.5	Structure of the Thesis	30
2	Problem Definition	32
2.1	Introduction	32
2.2	Assumptions and Notations	32
2.3	Problem Formulation	34
2.4	Approximation of Links as M/D/1 Queues and Its Limitations	34
2.5	Conclusions	36

3	Related Work	37
3.1	Introduction	37
3.2	Modelling Delay-Sensitive Traffic	38
3.2.1	Representation as Non-Concave Utility Functions	38
3.2.2	Representation as Hard Constraints	43
3.2.3	Comparison of the Modelling Approaches	45
3.3	Joint Rate Control and Scheduling for Elastic Traffic	46
3.3.1	Scheduling Solution Approaches	47
3.4	Joint Rate Control and Scheduling for Delay-Sensitive Traffic	49
3.4.1	Minimising Delay Using Virtual Data Rates	49
3.4.2	Minimising Network Congestion	49
3.4.3	Minimising Total Distortion for Video Transmissions	50
3.4.4	Providing Bounded Delay for Traffic with Elastic Bandwidth Requirements	51
3.5	Rate Control for Heterogeneous Traffic	52
3.5.1	Maximising Utility as a Function of Rate and Delay - The Concave Case	52
3.5.2	Maximising Utility as a Function of Rate and Delay - The Non-Concave Case	55
3.6	Conclusions	56
4	Alternative Problem Formulation	58
4.1	Introduction	58
4.2	The Alternative Optimisation Problem	58
4.3	Representation as a Scheduling Problem	60
4.3.1	Solution of the Multipath Rate Control Subproblem	61
4.3.2	Solution of the Scheduling Problem	63
4.4	Conclusions	67
5	Proposed Solution for Providing Bounded Delay	69
5.1	Introduction	69
5.2	Effect of Sources' Weights on Delay	70

5.3	Delay Regulation via Dynamic Adjustment of Sources' Weights . . .	72
5.4	Conclusions	77
6	Simulation Results	79
6.1	Introduction	79
6.2	Network Model	80
6.3	Implementation Using SimEvents	81
6.4	Results	85
6.5	Conclusions	89
7	Conclusions	99
7.1	Main Findings	99
7.2	Contributions to Knowledge	102
7.3	Limitations of the Work	103
7.4	Future Work	104
A	SimEvent Simulation Models	106
B	Mathematical Background	110
B.1	Sensitivity Analysis in Nonlinear Programming	110
B.2	Discontinuous Control	113
B.3	Matrix Analysis	114

List of Figures

1.1	An infrastructure-based wireless network (left) and a multi-hop wireless network (right)	11
1.2	Five-layer architecture for modular design of wireless networks [16] .	13
1.3	Transmission using node B as relay requires less power than direct transmission between nodes A and C	17
1.4	Interference region for the transmission via the relay nodes $R_1 \dots R_k$ is less than direct transmission between nodes A and C	18
1.5	Interdependencies between the functions within the layered architecture in multi-hop wireless networks	22
3.1	Utility functions for common classes of applications [34]	39
6.1	Network topology and alternative paths for source-destination pairs $A \rightarrow C$ and $E \rightarrow D$	82
6.2	Network contention graph and its maximal cliques	82
6.3	Path transmission rates when source weights are fixed at $w_1 = 2$ and $w_2 = 1$ and algorithms (4.13)-(4.14) are simulated	90
6.4	Link data rates when source weights are fixed at $w_1 = 2$ and $w_2 = 1$ and algorithms (4.13)-(4.14) are simulated	91
6.5	Path prices and normalised end-to-end delays when source weights are fixed at $w_1 = 2$ and $w_2 = 1$ and algorithms (4.13)-(4.14) are simulated	92
6.6	Packet end-to-end delay when algorithms (4.13)-(4.14) are simulated jointly with (5.5) with delay bounds $d_1 = d_2 = 1 \times 10^3$ msec .	93

6.7	Path transmission rates when (4.13)-(4.14) are simulated jointly with (5.5) with delay bounds $d_1 = d_2 = 1 \times 10^3$ msec, and flow $E \rightarrow D$ starts at time 2.5×10^4	94
6.8	Link data rates when (4.13)-(4.14) are simulated jointly with (5.5) with delay bounds $d_1 = d_2 = 1 \times 10^3$ msec, and flow $E \rightarrow D$ starts at time 2.5×10^4	95
6.9	Path prices when (4.13)-(4.14) are simulated jointly with (5.5) with delay bounds $d_1 = d_2 = 1 \times 10^3$ msec, and flow $E \rightarrow D$ starts at time 2.5×10^4	96
6.10	Packet end-to-end delays when (4.13)-(4.14) are simulated jointly with (5.5) with delay bounds $d_1 = d_2 = 1 \times 10^3$ msec, and flow $E \rightarrow D$ starts at time 2.5×10^4	97
6.11	Source weights when (4.13)-(4.14) are simulated jointly with (5.5) with delay bounds $d_1 = d_2 = 1 \times 10^3$ msec, and flow $E \rightarrow D$ starts at time 2.5×10^4	98
A.1	Simulation model of the network in Figure 6.1	107
A.2	Source 1 Rate Control subsystem	108
A.3	Link Price Update subsystem	109

List of Symbols

Symbol	Meaning
S	Set of traffic sources and its cardinality
s	Index for traffic sources
L	Set of links and its cardinality
l	Index for links
I_s	Set of alternative paths to the destination of source s and its cardinality
I	Total number of paths
i	Index for paths
R_i^s	$L \times 1$ vector defining the set of links used by path $i \in I_s$
$R_{l,i}^s$	Elements of vector R_i^s
R^s	$L \times I_s$ routing matrix for source s defined by $R^s = [R_1^s \dots R_{I_s}^s]$
R	$L \times I$ routing matrix for the network defined by $R = [R^1 \dots R^S]$
\mathbf{p}	$L \times 1$ vector of power assignments
p_l	Power assignment at link l
\mathbf{c}	$L \times 1$ vector of link data rates, or schedules
c_l	Data rate of link l
u	Rate-power function of the system which maps \mathbf{p} to \mathbf{c}
Π	Set of feasible power assignments
C	Set of feasible link data rates, or schedules
C^*	Set of optimal link data rates, or schedules
Co	Convex hull
x_i^s	Data transmission rate on path $i \in I_s$
x_s	Aggregate data transmission rate of source s
\mathbf{x}	$I \times 1$ vector of path transmission rates
y_l	Total traffic rate on link l
f_s	Utility function of source s

Symbol	Meaning
$\boldsymbol{\theta}$	$L \times 1$ vector of average queueing delay at links
θ_l	Average queueing delay experienced by a packet at link l
d_s	Upper limit on average end-to-end delay for source s
\mathbf{d}	$I \times 1$ vector with elements $d_i^s = d_s$, for all $i \in I_s$
w_s	Utility weight coefficient for source s
$\boldsymbol{\lambda}$	$L \times 1$ vector of link prices
λ_l	Price, or Lagrange multiplier associated with routing constraint at link l
\mathbf{q}	$I \times 1$ vector of path prices
q_i^s	Price of path $i \in I_s$
q_i^{s*}	Optimal price of path $i \in I_s$
q_s^*	Optimal path price for source s (i.e. q_i^{s*} are equal $\forall i \in I_s, x_i^{s*} > 0$)
$\boldsymbol{\mu}$	$I \times 1$ vector of Lagrange multipliers μ_i^s
μ_i^s	Lagrange multiplier associated with non-negativity of rate on path $i \in I_s$
$\mathbf{x}(\mathbf{c})$	Primal optimal solutions given \mathbf{c}
$(\boldsymbol{\lambda}(\mathbf{c}), \boldsymbol{\mu}(\mathbf{c}))$	Dual optimal solutions given \mathbf{c}
$(\mathbf{x}^*, \mathbf{c}^*)$	Primal optimal solutions
$(\boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$	Dual optimal solutions
$\Lambda(\mathbf{c})$	Set of optimal Lagrange multipliers associated with routing, given \mathbf{c}
$\mathbf{x}(\boldsymbol{\lambda})$	Path transmission rates given $\boldsymbol{\lambda}$
I_d	Set of disjoint paths and their cardinality
S_d	Set of sources with only disjoint paths and their cardinality
β	Step size parameter of link price update algorithm
γ	Step size parameter of the scheduling algorithm
α	Step size parameter of the delay regulator
∂_C	Clarke gradient
$\dot{\bar{V}}, \dot{\bar{\mathbf{q}}}$	Nonpathological derivatives of the maps V and \mathbf{q}
I_L	$L \times L$ Identity matrix
$\mathbf{x}^*(\mathbf{w})$	Primal optimal solutions given parameter \mathbf{w}
$\mathbf{q}^*(\mathbf{w})$	Optimal path prices given parameter \mathbf{w}
\mathbf{w}^*	$S \times 1$ vector of utility weight coefficients that guarantee bounded delay
$q_s(\mathbf{c}, \mathbf{w})$	Optimal path price for source s given \mathbf{c} and \mathbf{w}
$\mathbf{q}_w(\mathbf{c})$	Optimal path prices given \mathbf{c} and parameter \mathbf{w}

List of Abbreviations

Abbreviation	Meaning
MIMO	Multiple input multiple output system
UAV	Unmanned aerial vehicle
SINR	Signal to interference and noise ratio
BER	Bit error rate
PER	Packet error rate
TDMA	Time division multiple access
FDMA	Frequency division multiple access
CDMA	Code division multiple access
FH-CDMA	Frequency-hopping code division multiple access
CSMA	Carrier sense multiple access
AWGN	Additive white Gaussian noise
ARQ	Automatic repeat request
AODV	Ad hoc on-demand vector routing
DSR	Dynamic source routing
MDC	Multiple description coding
QoS	Quality of service
NUM	Network utility maximisation
MAC	Medium access protocol
FIFO	First in first out
TCP	Transmission control protocol
KKT	Karush-Kuhn-Tucker
VoIP	Voice over IP

Chapter 1

Introduction

1.1 Multi-hop Wireless Networks

A Multi-hop wireless network is composed of a cluster of wireless mobile nodes that form a network without a fixed infrastructure, using distributed control algorithms. A distinctive feature of such networks is multi-hop routing where any node can be a relay for the traffic of any other node in order to provide enhanced network coverage. Moreover, through multi-hop routing, data transmission between source and destination nodes is carried out via low power communication links between intermediate relay nodes. This leads to improved power efficiency – due to path loss exponential increase with distance – as well as overall network capacity – due to the reduced level of interference – as will be described in Section 1.2.3. Additionally, communication between sources and destinations can be carried out via alternative paths which can increase data transmission rates as well as robustness to the network topology changes (Figure 1.1). However, the lack of infrastructure in multi-hop wireless networks leads to more design complexities and less efficient utilisation of network resources, compared to the infrastructure-based wireless networks. Infrastructure-based wireless networks consist of stationary base stations positioned across a geographical area to optimise network coverage. Base stations are interconnected with high speed communication links, and connected to a backbone wired network. Each mobile node communicates directly with typically one

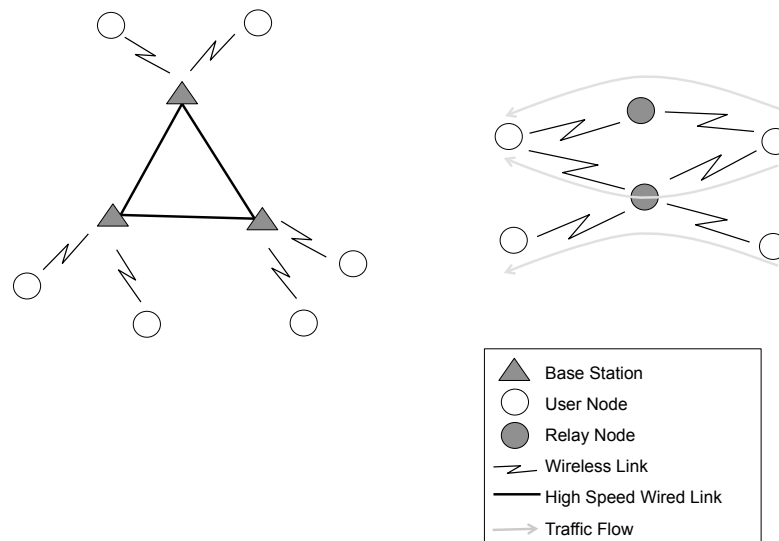


Figure 1.1: An infrastructure-based wireless network (left) and a multi-hop wireless network (right)

base station at a time via single-hop routes (Figure 1.1). Base stations provide access for mobile nodes to the network and perform all networking and control functions including transmission scheduling, power control, routing and handoff, in a centralised fashion. As a result, they lead to more efficient utilisation of network resources. Furthermore, as all computationally intensive functions are carried at base stations, the required computational tasks at mobile nodes are minimal. Finally, single-hop routes result in lower delay and loss, as well as higher data rates. In multi-hop wireless networks, however, networking and control functions are distributed among all wireless nodes, any node can be a relay for other nodes traffic, and nodes are typically in motion.

Despite its performance limitations, multi-hop wireless networking is the only viable technology in wireless communication applications where low-cost, rapid deployment and configuration, and robustness are critical factors. An example

relevant to the subject of this thesis is providing communication support for distributed control systems with remote plants, sensors and actuators, in particular, coordinated control of unmanned mobile units such as unmanned aerial vehicles (UAVs) and vehicular networks in automated highway systems. Such systems require that sensor and actuator signals to be delivered to the controller with a small and short delay, but can adapt to various data rates [16]. Other application example is real-time interactive audiovisual communication in disaster relief or military applications. Similar to the distributed control systems these applications impose strict limit for maximum packet end-to-end delay. In addition, they have high but flexible data rate requirements [24, 41, 37]. Supporting applications with high data rates and bounded delay requirements in multi-hop wireless networks is highly challenging due to the performance limitations described earlier and requires a joint optimisation design approach, as will become evident in the following sections.

1.2 Multi-hop Wireless Networks Design

The main design issues in multi-hop wireless networks can be best described by dividing the main network functions into the conventional five-layer architecture used for modular design of wireless networks, as described in [16] (Figure 1.2).

1.2.1 Physical Layer

The physical layer involves functions that deal with transmitting bits over a point-to-point wireless link. They include modulation/detection, coding, power control and multiple input multiple output (MIMO) systems. Digital modulation consists of mapping the information bits into an analog signal for transmission over the channel. Detection consists of determining the original bit sequence from the received signal. Digital modulation techniques are chosen based on the performance characteristics including data rate, spectral efficiency, power efficiency and robustness to channel impairments. Coding enables bit errors to be either detected or corrected by a decoder in the receiver. The resulting performance enhancement,

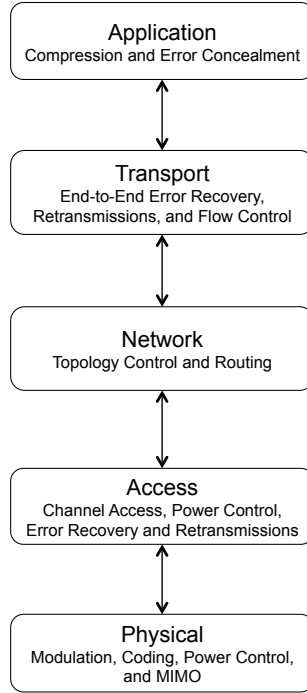


Figure 1.2: Five-layer architecture for modular design of wireless networks [16]

however, leads to increased complexity, decreased data rate or increase in signal bandwidth. Power control adjusts the transmit power of all nodes such that signal to interference and noise ratio (SINR) at each receiving node does not fall below its minimum required level for acceptable performance. Power control is also an access layer functionality and will be described in more detail as part of access layer function. MIMO refers to systems with multiple antennas at the transmitter and the receiver. Multiple antennas can be used to increase link data rates through multiplexing, providing diversity to fading to reduce the average bit error rates (BER), and providing directionality to reduce fading and interference to other signals.

The design choices at the physical layer impact the other layers in several ways. For example, they determine the link packet error rate (PER) which affect the retransmission at the access layer. In addition, multiple antennas increase the

link data rate and reduce interference to other links which impacts all the protocol layers. Similarly, the transmission power together with adaptive modulation and coding for a node determines the collection of nodes it can reach in a single hop and therefore affects the context in which the higher layers operate.

1.2.2 Access Layer

The access layer controls the allocation the available spectrum to the users and ensures successful reception of packets over this spectrum. Signalling dimensions are allocated using either multiple access or random access techniques. In multiple access signalling dimensions are divided into dedicated channels which are assigned to different users. This includes orthogonal division along the time axis, as in time division multiple access (TDMA), orthogonal division along the frequency axis, as in frequency division multiple access (FDMA), and orthogonal or non-orthogonal division along the code axis, as in code division multiple access (CDMA). In random access channels are allocated to the active users dynamically, as in ALOHA, CSMA, and scheduling. Multiple access methods are suitable for applications with continuous data transmissions and delay constraints, while random access methods are more suitable for users with bursty data transmissions. The access layer also involves control mechanisms for channel assignment to users and their admission into the system.

Power control is also a function of the access layer. As mentioned previously, power control ensures link required SINR levels are satisfied by adjusting the transmission power of all transmitting nodes. Link required SINR levels are determined by link performance requirements such as connectivity in topology control (a function of the network layer), and link data rates decisions. Precisely, let L be the number of links in a multi-hop wireless network. The SINR at link l is then given by

$$\gamma_l = \frac{g_{ll}p_l}{n_l + \rho \sum_{j \neq l} g_{lj}p_j}$$

where $g_{lj} > 0$ is the channel power gain from the transmitter of the link j to the receiver of the link l , n_l is noise power at receiver on the link l , and ρ is the

interference reduction due to signal processing. For example, $\rho \approx \frac{1}{G}$ for CDMA with processing gain G , and $\rho = 1$ for TDMA. The data rate of link l can be modelled as function of its SINR, γ_l , as well as physical layer parameters. For example, the Shannon capacity of link l for additive white Gaussian noise (AWGN) channel is given by $c_l = W \log_2(1 + \gamma_l)$, where W is the channel bandwidth. Or in [40], the data rate of link l is assumed to be given by $c_l = W \log_2(1 + \frac{\gamma_l}{\Gamma})$, where Γ is a parameter determined by the physical layer design.

Let γ_l^* be the minimum required SINR for link l . The SINR constraints are then given by $(I - F)\mathbf{p} \geq \mathbf{u}$, where \mathbf{p} is the $L \times 1$ vector of transmission powers with elements p_l , \mathbf{u} is an $L \times 1$ vector with elements $\frac{\gamma_l^* n_l}{g_{ll}}$, and F is an $L \times L$ matrix with elements

$$F_{lj}^s = \begin{cases} 0 & l = j, \\ \frac{\gamma_l^* n_l \rho}{g_{ll}} & l \neq j. \end{cases}$$

If the eigenvalues of F are strictly inside the unit circle then there exists a power assignment \mathbf{p} that supports the required SINR for all links, and $\mathbf{p}^* = (I - F)^{-1}\mathbf{u}$ is the minimal (Pareto optimal) solution. In [15], a simple distributed power control algorithm is proposed which converges to the minimal power solution \mathbf{p}^* when feasible solution exists. This algorithm is extended in [3] to incorporate admission control. However, it may be inefficient to activate all links concurrently at power levels \mathbf{p}^* , and combining power control and transmission scheduling can improve power efficiency while satisfying minimum required SINR levels, as discussed in [12]. For example in [12] an optimal transmission scheduling and power control strategy is proposed which minimises total average transmission power while ensuring the minimum required link data rates. Or in [13], a joint scheduling and power control strategy is proposed which maximises network throughput and reduces power consumption.

Retransmission of corrupted packets, which is referred to as ARQ protocol, is also a function of the access layer. In this protocol, using the error detection code the receiver determines if there are corrupted bits in the packet that cannot be corrected. The receiver then discards the corrupted packet and sends a retransmission request to the transmitter. Alternatively, to improve the network throughput,

the receiver can save the original packet and use a form of diversity to combine it with the retransmitted packet, or the incremental redundancy technique can be used where instead of retransmission of the entire packet, the transmitter just sends some additional coded bits to provide a stronger error correction capability for the original corrupted packet.

1.2.3 Network Layer

The role of the network layer is to establish and maintain end-to-end connections in the network. The associated functions include topology control, routing and dynamic resource allocation. Topology control constitutes mechanisms for adjusting and coordinating nodes transmission powers in order to generate networks with the desired properties such as connectivity, while reducing overall power consumption and/or increasing network capacity [33]. Roughly speaking, a network is connected if for any two nodes in the network there exists a path through which they can communicate. The transmission power determines the transmission range of a node within which direct communication with other nodes is possible given minimum link performance requirements such as link data rates and BER. Thus, in the process of topology control each node identifies the network nodes that it can directly communicate with.

The impact of topology control on improving network energy efficiency and capacity can be demonstrated by means of the example shown in Figures 1.3 and 1.4 [33]. Figure 1.3 shows alternative transmission routes between nodes A and C , where d_{ij} denotes the distance between nodes i and j . It is assumed that both nodes B and C are within the transmission range of node A . Therefore, there are two alternative paths between nodes A and C : the multi-hop path $A \rightarrow B \rightarrow C$ using node B as relay, and the direct path $A \rightarrow C$. Assuming the free space propagation model; i.e. $P_r \propto \frac{P_t}{d^2}$ where P_r and P_t are the received and transmitted power, respectively, and d is the distance between the transmitter and the receiver; the power required for successful transmission between two nodes is proportional

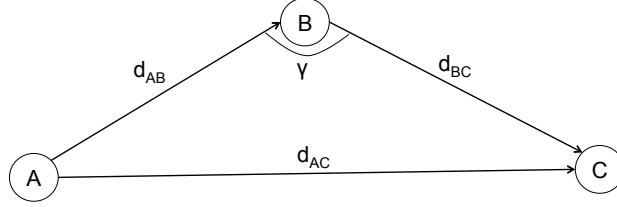


Figure 1.3: Transmission using node B as relay requires less power than direct transmission between nodes A and C

to the square of their distances. From Figure 1.3 it follows that

$$\begin{aligned} d_{AC}^2 &= d_{AB}^2 + d_{BC}^2 - 2d_{AB}d_{BC} \cos \gamma \\ &> d_{AB}^2 + d_{BC}^2 \end{aligned}$$

Thus, transmission using node B as relay requires less power than direct transmission between nodes A and C .

Similarly, in the network shown in Figure 1.4 there are two alternative transmission strategies between nodes A and C : a direct transmission from A to C , and a multi-hop transmission using nodes $R_1 \dots R_k$ as relays. To compare the impact of the two alternative transmission strategies on the network capacity, the notion of interference region based on the protocol interference model introduced in [17] is used. According to the protocol interference model, the interference region for a receiving node j is defined as a circle of radius $(1 + \eta)d_{ij}$ centred at node j , where $\eta > 0$ is a constant that depends on the features of the wireless transceiver and d_{ij} is the distance between the transmitting node i and receiving node j . It specifies the region where no other node can transmit simultaneously in order for transmission from node i to node j to be successful. The interference region measures the

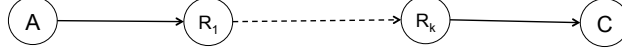


Figure 1.4: Interference region for the transmission via the relay nodes $R_1 \dots R_k$ is less than direct transmission between nodes A and C

amount of wireless medium used by a certain transmission and thus can be used as a measure of the network capacity. Using the above model, the interference region for direct communication from node A to node C is given by $\pi d_{AC}^2 (1 + \eta)^2$. On the other hand, assuming $d_{AR_1} = d_{R_1 R_2} = \dots = d_{R_k C} = \frac{d_{AC}}{k+1}$, the interference region for the transmission via the relay nodes $R_1 \dots R_k$ is given by $(k+1)\pi \left(\frac{d_{AC}}{k+1}\right)^2 (1 + \eta)^2$ which, by Holder's inequality, is less than $\pi d_{AC}^2 (1 + \eta)^2$. This implies that using multi-hop transmissions can reduce total interference range and thus increase the overall network capacity.

After performing topology control, the power control algorithms described in Section 1.2.2 can be activated to achieve link required SINR levels with high power efficiency. However, in the presence of node mobility, nodes transmission power as well as physical layer parameters have to be dynamically adjusted to compensate for the variations in link gains and hence maintain connectivity and the link required SINR levels.

Routing determines how packets are routed through the network from their source to their destination. Routing protocols for multi-hop wireless networks can be classified as flooding, proactive (centralised, source-driven, and distributed),

and reactive. Flooding is based on broadcasting a packet to all neighbouring nodes which then broadcast the packet to their neighbouring nodes. This process repeats until the packet reaches its destination. This approach requires little communication overhead and thus is suitable for highly mobile networks where network topology changes rapidly. However, flooding is highly inefficient in terms of bandwidth utilisation and energy consumption and as a result is only practical in small networks.

In centralised routing approach, information about network topology and channel condition are determined by each node and forwarded to a central node, which then based on an optimisation criterion computes the routing tables for all nodes and communicates them to the nodes. Centralised approach results in globally optimal routes, however, it leads to significant communication overhead for exchanging routing information and as a result it cannot adapt to the rapid changes in network topology or channel conditions. For this reason, this approach is only suitable for small networks. In source-driven routing, each node receives the network topology information and computes optimal routes to its destination. This method also involves periodic exchange of routing information, which leads to high communication overhead. In distributed routing, each node communicates its connectivity information to its neighbouring nodes, and each node determines its local routing information, i.e. the next hop in the route, using this local information. This approach requires low communication overhead and as such can adapt quickly to link and connectivity changes. On the other hand, the computed routes are suboptimal in general.

Both centralised and distributed approaches are based on maintaining up-to-date routing information. In reactive (on-demand) routing approaches, however, routes computation is initiated by a source node that has traffic to send, via a route discovery process. The process is completed once a route or all possible routes are discovered. The routes are then maintained by a route maintenance procedure until they are no longer needed, or destination becomes inaccessible. This approach leads to globally optimal routes with little overhead, since routes are maintained when only in use, but involves significant initialisation delay. Examples of reactive

routing include ad hoc on-demand vector routing (AODV) and dynamic source routing (DSR).

As will be explained in Section 1.2.6, due to high coupling between the layer functions, computation of optimal routes has to be carried out in conjunction with design choices at other layers, in order to optimise overall network performance. This calls for maintaining multiple paths between source-destination pairs which enables splitting the traffic load among multiple paths to optimise network-wide performance objectives. Multi-path routing can also be used to increase the probability of packet reception in networks with rapidly changing topologies, by transmitting duplicates of a packet over multiple paths between the source and the destination [38]. Multi-path routing in conjunction with network coding has also been shown to improve multi-hop wireless network throughput, when multiple multicast sessions are present [25]. Network coding is based on fusion of data received from multiple routes which can then be decoded by intermediate or destination nodes, as necessary.

1.2.4 Transport Layer

The transport layer consists of end-to-end functions of error recovery, retransmission, reordering and flow control. The error recovery and retransmission involve mechanisms that check for corrupt or lost packets on the end-to-end route and send a request for retransmission to the source node. The reordering function is responsible for ordering packets that arrive out of order due to multi-path routing, delay, congestion, packet loss, or retransmission; before being passed to the application layer. Flow control allocates the flow associated with the application layer to different routes, based on an optimisation criterion.

1.2.5 Application Layer

The application layer involves functions that deal with generating the data to be transmitted, and processing the corresponding data received over the network. These functions include data compression, error correction and concealment. The

level of data compression depends on the application's selected trade-off between data rate and robustness to the changes in network conditions. High level of compression reduces the required data rate but increases the data sensitivity to error since it removes most of the redundancy. Applications such as voice and video can tolerate some errors without significant degradation in their perceived quality. In contrast, data applications cannot tolerate any packet loss and as a result any lost or corrupted packet has to be retransmitted.

The application layer can also perform multiple description coding (MDC) whereby multiple description of data are generated and the original data can be reconstructed from any description with some loss. Using multi-path routing, each description can then be sent over a different path to provide diversity in the presence of network performance degradation.

Applications typically have performance requirements for end-to-end data rates and delay, which are referred to as quality of service (QoS). Examples include distributed control applications [16], which require bounded end-to-end delay but may be able to tolerate a lower data rate via a coarser quantisation of sensor data. Similarly, interactive real-time voice and video applications require bounded end-to-end delay but can adapt to various data rates using various encoding quantisations [24, 41, 37].

1.2.6 Motivations for Cross-Layer Design

The functions within the layered architecture described in the previous sections are tightly coupled in multi-hop wireless networks, as shown in Figure 1.5. The network performance characteristics, e.g. end-to-end delay and data transmission rates, are dependent on the network congestion level, which is determined by rate of the traffic flow entering the network, the distribution of load across the network links and the network link capacities. The flow rate is regulated by the flow control function at the transport layer, given the QoS requirements from the application layer and current network performance characteristics such as a measure of congestion level. Distribution of traffic load is performed at the network layer by the multi-path routing algorithm, based on the optimisation criterion, e.g.

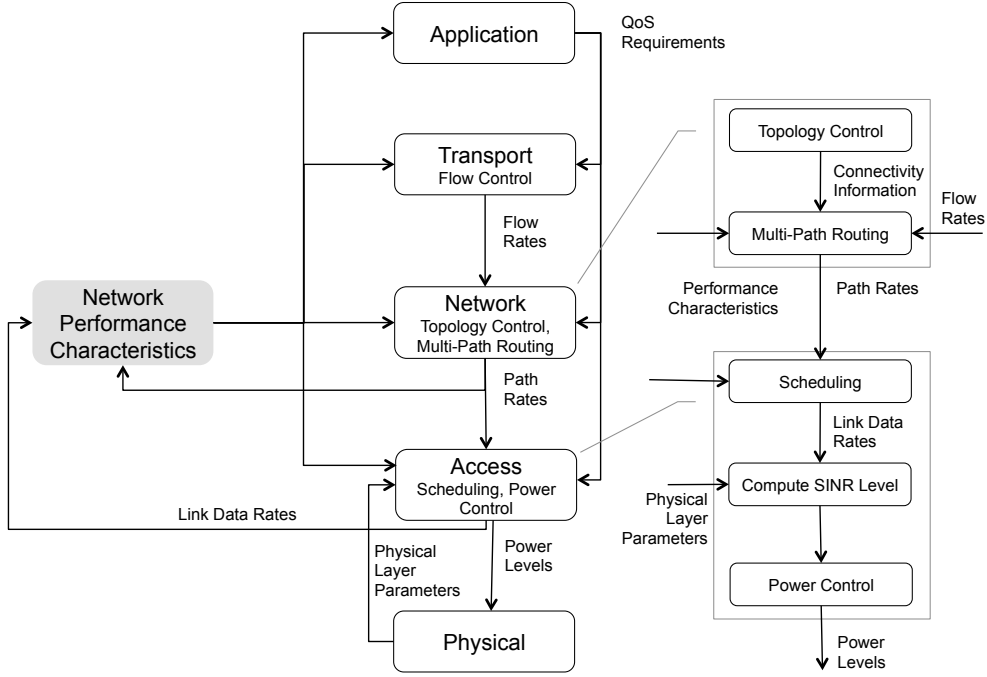


Figure 1.5: Interdependencies between the functions within the layered architecture in multi-hop wireless networks

minimum network congestion, and current network performance characteristics. The performance of the multi-path routing algorithm is dependent on the network connectivity which is determined by the topology control algorithm. Given the current network performance such as link congestions levels, link capacities (data rates) are updated by the scheduling algorithm based on the optimisation criterion. Given the physical layer parameters, the updated link data rates can be attained at certain link SINR thresholds using power control algorithms. The transmission powers can be jointly optimised with other physical layer parameters, e.g. modulation and coding, to achieve the required link data rates. Finally, the current performance characteristics can also be used at the application layer to determine parameters such as data compression rates, encoding quantisation rate, or to decide whether to transmit data over multiple paths via MDC mechanism.

The high level of interdependencies between the functions in the layered architecture motivates the joint optimisation of the design variables approach to solve a network design problem. Specifically, in this approach the network design problem is represented as a network utility maximisation (NUM) problem [8] where the end-user application preferences are captured as the objective function, design constraints as the constraint set, and design freedoms as the optimisation variables. Using the theory of decomposition for non-linear optimisation, the network design problem can then be systematically decomposed into various functional modules alternatives that correspond to a layered architecture. Each functional module can be further decomposed into distributed computation and control over distant network elements. Moreover, optimal solutions serve as benchmark for heuristic layered architectures.

This research is based on the following NUM formulation presented in [29]

$$\begin{aligned}
 & \max_{\mathbf{x}, \mathbf{c}} \quad \sum_{s \in S} f_s(x_s) \\
 & \text{subject to} \quad R\mathbf{x} \leq \mathbf{c} \\
 & \quad \mathbf{c} \in \text{Co}(C) \\
 & \quad \mathbf{x} \geq 0
 \end{aligned}$$

The formal definition of the notations are provided in Chapter 2. Briefly, x_s denote data transmission rate for the traffic source $s \in S$, where S is the set of traffic sources. The utility function f_s represent the elastic QoS requirement for the traffic source s . Here, the utility function is assumed to be a measure of the optimality of resource allocation efficiency. In other words, it defines the QoS value a traffic source attributes to a particular transmission rate. The utility function is typically assumed in the congestion control literature [35, 29] to be smooth, increasing, and strictly concave function of the transmission rate. This shape of utility function has been shown to lead to optimal resource allocation that can achieve various fairness objectives [35, 29]. Furthermore, as discussed in [34], this shape of utility function is typical (except for very small bandwidths) of the utility (performance) of rate adaptive real-time applications which can adjust their transmission rate by applying different encoding levels. With the exception

of very small bandwidths, the marginal rate of utility (performance) enhancement is diminishing as bandwidth increases. At high bandwidths the marginal utility of additional bandwidth is very small since the signal quality is higher than what human can perceive. Examples of such applications include voice over IP (VoIP) application as described in [39]. Normally, acceptable voice quality requires a transmission rate of about 64 kb/s to transmit an 8 kHz 8 bit signal. While it is possible to cope with a little less quality than this, the utility of transmitting voice at less than 64 kb/s drops sharply as the transmission rate is decreased. Also, the quality of voice does not increase much if the transmission rate is increased beyond 64 kb/s, so the utility function remains almost flat above 64 kb/s.

R denotes the routing matrix that defines each source's available paths to its destination. It is assumed that routing tables are already computed by a source-driven routing algorithm, based on network connectivity resulted from a topology control algorithm (Section 1.2.3). Thus, the key assumption in the above formulation is that variations at the physical layer are such that they can be compensated at the physical layer without affecting the network topology within the timescale of the above problem. \mathbf{x} denotes path transmission rates. \mathbf{c} denotes the link data rates and assumed to be a function of global power assignments, or any other resource control decisions such as link activation and inactivation, or retransmission probability in random access MAC protocols, generally denoted by \mathbf{p} . This implies $\mathbf{c} = u(\mathbf{p})$. The function u captures the cross-layer control decisions at both physical and access layers (Sections 1.2.1 and 1.2.2). Specifically, by choosing appropriate physical layer parameters, e.g. modulation and coding, link data rates are mapped to link SINR levels and the corresponding power assignments via u^{-1} . Let Π be the set of feasible power assignments, then $C = \{u(\mathbf{p}), \mathbf{p} \in \Pi\}$ is the set of feasible link data rates, or schedules. Link data rates are constrained by the convex hull of C denoted by $\text{Co}(C)$, which captures the time interleaving of the feasible schedules. It is assumed that variations in link SINR levels can be compensated quickly by power control or adjusting other physical layer parameters without affecting the link data rates within the timescale of the above problem. In addition to the constraints related to the network structure, per-user inelastic QoS

constraints can be added to the above formulation. The optimisation variables, or the design freedoms, in the above NUM problem are the path transmission rates \mathbf{x} , and link data rates \mathbf{c} . The utility functions f_s , $s \in S$, the routing matrix R , and the set of feasible schedules C are treated as constants over the problem timescale.

1.3 Research Objectives

The main focus of this research is the problem of supporting traffic with high data rate requirements and hard end-to-end delay constraints in multi-hop wireless networks, using source data transmission rates and link data rates as the key design variables. Before presenting the research objectives, first, different load control mechanisms for QoS provision are described and the preferred approach is introduced. The design problem and its underlying assumptions are then described followed by a summary of the limitations of the current solutions.

1.3.1 Background on QoS-Oriented Load Control

As discussed in [39], based on the type of QoS guarantees that can be provided, load control mechanisms can be classified into Bandwidth reservation, and best-effort schemes. Bandwidth reservation based schemes are suitable for supporting traffic with minimum bandwidth requirements. In these schemes the network is required to maintain information about each traffic flow and decide whether a traffic flow can be admitted so that all the admitted traffic flows can be guaranteed their minimum required bandwidth. Inevitably, the rejection of some of the traffic results in inefficient utilisation of network capacity. More importantly, Bandwidth reservation based schemes lead to a significant communication overhead in networks with dynamic settings and as a result are unsuitable for multi-hop wireless networks due to their dynamic and distributed characteristics. In contrast, best-effort schemes are suitable for supporting traffic whose demand for bandwidth is elastic, but their perceived QoS referred to as their utilities, are generally assumed to be an increasing and concave function of their transmission rates. As discussed in Section 3.2.1 and [34], for elastic traffic the overall network QoS is maximised

by admitting all traffic flows and allocating network capacity based on the traffics' perceived QoS. This is the main design principle of best-effort schemes. The key advantage of the best-effort approach is that it enables QoS optimal allocation of network resources using simple and distributed algorithms [39], and as a result is a suitable approach for QoS provision in multi-hop wireless networks. The basic NUM formulation presented in Section 1.2.6 results in a best-effort solution and hence forms the basis of this research.

1.3.2 Problem Description and Assumptions

The fundamental assumption in this research is that all incoming traffic have elastic bandwidth requirements, but their perceived QoS due to bandwidth, i.e. excluding the effect of end-to-end delay, are increasing and concave functions of their transmission rates. Furthermore, all or some of the incoming traffic impose a strict limit on the average end-to-end queueing delay. These requirements are typical of QoS requirements of delay-sensitive applications such as distributed control systems and real-time interactive audiovisual communication, as explained in Section 1.1.

As explained in Section 1.2.6, it is assumed that routing tables are already computed by a source-driven routing algorithm, and remain unchanged within the time horizon of the problem. Furthermore, it is assumed the set of feasible link data rates, or schedules, are known and remain constant over the time horizon of the problem.

The key design variables are assumed to be the path transmission rates and link data rates, or schedules. The main research objective is then to develop a path rate control and scheduling strategy that ensures average end-to-end queueing delays do not exceed their imposed upper bounds and maximises the aggregate utility, or the perceived QoS due to bandwidth, of all incoming traffic. Moreover, the strategy should enable distributed implementation with low communication overhead in order to be deployable in a multi-hop wireless network setting.

1.3.3 Limitations of Current Solutions

As will be described in Section 3, the network utility maximisation based approaches to support delay-sensitive traffic have been predominantly based on either reducing link utilisation, or approximation of links as $M/D/1$ queues. The former approach normally leads to nearly zero queue lengths in the long term due to reduced link utilisation, but provides no control over the transient behaviour of packet delays. The key assumptions behind the latter approach are unrealistic since the traffic at entry points are regulated by the rate controller and are deterministic; multi-hop networks are composed of mostly disjoint paths which comprise serial links, and the traffic entering the links can be further regulated to limit its burstiness. Moreover, the algorithms for rate control and scheduling are typically implemented as feedback control systems where path rates and link data rates are regulated based on the current congestion level at links. The delay caused by the burstiness of the arriving traffic at each link can therefore be assumed to be negligible and consequently the average delay a packet experiences at equilibrium is primarily a function of number of packets in the system at equilibrium, which is determined by the dynamics of the rate control and scheduling algorithms at their transient state. The other limitation of the $M/D/1$ approximation is that it results in under-utilised links under optimal resource allocation.

1.3.4 Research Objectives

Based on the problem description and the shortcomings of the current solutions described in Sections 1.3.2 and 1.3.3, respectively, the research assumptions and objectives are summarised as follows:

Assumptions

- network topology remain fixed over the time horizon of the problem,
- routing tables are already computed by each traffic source and remain unchanged over the time horizon of the problem,

- the set of feasible link data rates, or schedules, are known and remain constant over the time horizon of the problem,
- all incoming traffic have elastic bandwidth requirements, but their perceived QoS due to bandwidth are increasing and concave functions of their transmission rates,
- main design freedoms are the path transmission rates and link data rates, or schedules.

Objectives

- to develop a path rate control and scheduling strategy that
 - ensures average end-to-end queueing delays do not exceed their imposed upper bounds,
 - maximises the aggregate utility, or the perceived QoS due to bandwidth, of all incoming traffic,
 - enables distributed implementation with low communication overhead in order to be deployable in a multi-hop wireless network setting,
 - leads to maximal link capacity utilisation,
 - has controllable transient behaviour.

1.4 Summary of Contributions

Given the limitations of the conventional modelling of the delay constraints based on $M/D/1$ queue approximation of links, an alternative formulation to the original optimisation problem is considered where delay constraints is omitted and the source utility functions are multiplied by weight coefficients. The alternative optimisation problem is then transformed into a master scheduling problem and the well-known multi-path rate control with fixed link data rates subproblem.

The multi-path rate control subproblem is solved using a duality-based algorithm, which leads to equilibrium link queueing delays that are proportional to

the optimal dual variables, or link prices. While in general optimal path rates are not unique, and path rates in duality-based algorithms do not converge and continuously oscillate, conditions on the number of disjoint paths are derived that guarantee unique optimal path rates. The underlying theoretical conditions that guarantee unique optimal path rates can be further used for future work to design a multi-path rate control algorithm that converges the unique optimal path rates, given the conditions on the number of disjoint paths are guaranteed by the topology control algorithm.

A distributed algorithm for the master scheduling problem is then proposed, and is shown to converge to the optimal data rates. The proposed algorithm incorporates the solution of a well-known scheduling problem, for which efficient and distributed solutions have been developed in several cases.

For the alternative optimisation problem, derived bounds on the sensitivity of optimal path prices and aggregate source rates to the variations of utility weight coefficients indicate that optimal path prices for each source increase with the source's weight coefficient. Given the correlation between path queueing delays and path prices in the proposed scheduling algorithms, an alternative approach is then presented where utility weight coefficients are used as control variables to regulate end-to-end queueing delays in the scheduling algorithms. Specifically, an integral controller is incorporated in the scheduling algorithm whereby each source regulates the queueing delay on its paths at the desired level, using its weight coefficient as the control variable.

The conditions under which the proposed joint scheduling algorithm and delay regulator achieve asymptotic regulation of end-to-end delay are examined. The proposed joint scheduling algorithm and delay regulator meet the objectives stated in Section 1.3.4 since they regulate end-to-end queueing delay at the desired levels, can be implemented distributively, and lead to maximal link utilisation. For future work, linearisation and linear system design methods can further be used to adjust the delay regulator parameter for the desired transient behaviour.

Simulation experiments show that the presence of feedback error the proposed scheduling algorithm converges to the optimal solution of the alternative optimisa-

tion problem. In addition, the proposed joint scheduling algorithm and delay regulator achieve asymptotic regulation of end-to-end delay. Simulation experiments further demonstrate the deficiencies of other alternative approaches compared with the proposed approach.

1.5 Structure of the Thesis

The rest of this thesis is organised as follows. In Chapter 2, the problem is presented formally as a network utility maximisation problem, and the limitations of delay models based on approximation of links as $M/D/1$ queues is discussed. After describing the assumptions and notations in Section 2.2, in Section 2.3 the network utility maximisation formulation is presented. The limitations of $M/D/1$ approximation of links for delay estimation are discussed in Section 2.4.

In Chapter 3 the predominant approaches that aim to address the requirements of delay-sensitive traffic using the NUM framework, as well as their limitations, are described. In Section 3.2, the two approaches to model delay-sensitive traffic, namely, representation as non-concave utility functions and hard constraints specification, are explained. In Section 3.3, the well-known distributed solutions for the problem of joint rate control and scheduling for elastic traffic, are presented. Various approaches that tackle the same problem but for inelastic traffic, including minimising delay using virtual data rates, minimising network congestion, and minimising network distortion for video traffic, are introduced in Section 3.4. Finally, distributed rate control algorithms for networks with heterogeneous traffic, for cases where delay-sensitivity is modelled as a concave, as well as a non-concave function, are discussed in Section 3.5.

Given the limitations of the original optimisation problem discussed in Chapter 2, in Chapter 4 an alternative formulation and its proposed solution is presented, on which the proposed approach for providing bounded delay will be based. In Section 4.2 the proposed alternative optimisation problem, in which the delay constraints are omitted and utility functions are multiplied by weight coefficients, is presented and its properties are examined. In Section 4.3 the scheduling representation of

the alternative problem is introduced, which is based on vertical decomposition of the problem into master scheduling problem and the well-known multipath rate control with fixed link data rates subproblem. After deriving conditions for the uniqueness of the solutions of the multipath rate control subproblem, the proposed distributed algorithm for solving the scheduling problem is presented.

In Chapter 5, the proposed solution to the original optimisation problem is developed, with focus on the performance objectives described in Section 1.3.4. The proposed solution exploits the properties of the alternative optimisation problem and its proposed solution described in Chapter 4. In Section 5.2, bounds on the sensitivity of optimal path prices and aggregate source rates to the variations of utility weight coefficients in the alternative optimisation problem are derived. Based on the sensitivity results, in Section 5.3 an algorithm for providing bounded end-to-end queueing delays as well as other performance objectives is proposed, which is integrated in the scheduling algorithms developed in Chapter 4.

In Chapter 6, simulation experiments are performed to address three fundamental questions. Firstly, to illustrate that the proposed algorithm for solving the scheduling problem presented in Chapter 4 converge to its optimal solutions, despite using approximate values of link prices computed by the inner layer rate control algorithm. Secondly, to illustrate that the proposed joint scheduling algorithm and delay regulator in Chapter 5 can regulate packet end-to-end latency, using an estimation of end-to-end delay as feedback in the scheduling algorithm, and to compare their performance against the previously proposed main approaches to support delay-sensitive traffic. Finally, to assess the dynamic behaviour of the proposed algorithms when network configuration changes. In Section 6.2, the simulated network and its mathematical model are described. The SimEvents implementation of the proposed algorithms in Chapters 4 and 5 for the network model is described in Section 6.3. The result of the simulation experiments is presented in Section 6.4. The conclusions are presented in Section 6.5.

Chapter 2

Problem Definition

2.1 Introduction

In this chapter the problem is presented formally as a network utility maximisation problem, and the limitations of delay models based on approximation of links as $M/D/1$ queues is discussed. After describing the assumptions and notations in Section 2.2, in Section 2.3 the network utility maximisation formulation is presented. The limitations of $M/D/1$ approximation of links for delay estimation are discussed in Section 2.4.

2.2 Assumptions and Notations

Throughout the text, vectors are denoted by boldface lowercase letters, and matrices and sets by capital letters. For simplicity, the same notations are used to denote the sets and their cardinality.

This thesis considers the problem of rate control and scheduling for simultaneous transmissions of multiple delay-sensitive traffic over a multi-hop wireless network. Let S be the set of sources which generate the delay-sensitive traffic and L be the set links which constitute the multi-hop wireless network. Each source $s \in S$ has multiple alternative paths to its destination denoted by I_s . The set of

links used by each path $i \in I_s$ are defined by the $L \times 1$ vector R_i^s with elements

$$R_{l,i}^s = \begin{cases} 1 & \text{if path } i \in I_s \text{ uses link } l, \\ 0 & \text{otherwise.} \end{cases}$$

The $L \times I_s$ routing matrix for source s is subsequently defined by $R^s = [R_1^s \dots R_{I_s}^s]$, and the $L \times I$ routing matrix for the network, where $I = \sum_{s \in S} I_s$, by $R = [R^1 \dots R^S]$.

Let p_l be the power assignment, or any other resource control decisions such as activation and inactivation, and retransmission probability in random access MAC protocols, and c_l be the data rate at link l . Based on the NUM formulation presented in [29], link data rates are assumed to be a function of global power assignments, i.e. $\mathbf{c} = u(\mathbf{p})$. As explained in Section 1.2.6, the function u captures the cross-layer control decisions at both physical and access layers (Sections 1.2.1 and 1.2.2). Specifically, by choosing appropriate physical layer parameters, e.g. modulation and coding, link data rates are mapped to link SINR levels and the corresponding power assignments via u^{-1} . Let Π be the set of feasible power assignments, then $C = \{u(\mathbf{p}), \mathbf{p} \in \Pi\}$ is the set of feasible link data rates, or schedules. The convex hull of C denoted by $\text{Co}(C)$ captures the time interleaving of the feasible link data rates, and is assumed to be closed and bounded.

Let x_i^s be the data transmission rate on path $i \in I_s$, and $x_s = \sum_{i \in I_s} x_i^s$ be the aggregate data transmission rate of source s . It is assumed that each source $s \in S$ gains a utility $f_s(x_s)$ at rate x_s . f_s are assumed to be twice continuously differentiable, strictly concave, and increasing for all $s \in S$. Furthermore $f_s'' < 0$. As explained in detail in Section 1.2.6, this shape of utility function is the typical assumption in the congestion control literature since it leads to various fairness objectives at optimality, and furthermore represents the performance of rate adaptive real-time applications.

It is assumed that the average delay experienced by a packet at link l is given by $\theta_l(c_l, y_l)$, where $y_l = R_l \mathbf{x}$ is the total traffic rate on link l . Furthermore, $\theta_l(c_l, y_l)$ are differentiable, decreasing in c_l and increasing in y_l , for all $l \in L$. Let $\boldsymbol{\theta}$ be the $L \times 1$ vector with elements $\theta_l(c_l, y_l)$, for all $l \in L$. The average end-to-end delay

on path $i \in I_s$ of source $s \in S$ is then given by $R_i^{sT} \boldsymbol{\theta}$, and is assumed to be upper bounded by d_s . Let \mathbf{d} be an $I \times 1$ vector with elements $d_i^s = d_s$, for all $i \in I_s$.

2.3 Problem Formulation

The objective is to find data transmission rates \mathbf{x} and link data rates \mathbf{c} such that

$$\max_{\mathbf{x}, \mathbf{c}} \quad \sum_{s \in S} f_s(x_s) \quad (2.1)$$

$$\text{subject to} \quad R\mathbf{x} \leq \mathbf{c} \quad (2.2)$$

$$\mathbf{c} \in \text{Co}(C) \quad (2.3)$$

$$R^T \boldsymbol{\theta} \leq \mathbf{d} \quad (2.4)$$

$$\mathbf{x} \geq 0 \quad (2.5)$$

The optimisation objective (2.1) is to maximise the aggregate utility of all sources. Constraint (2.2) requires that the traffic rate entering each link not to exceed its allocated data rate. Constraint (2.3) restricts link data rates to the convex hull of feasible link data rates. Constraint (2.4) imposes an upper bound on the average end-to-end delay faced by a packet on individual paths.

2.4 Approximation of Links as M/D/1 Queues and Its Limitations

As will be discussed in Chapter 3, estimation of the average delay experienced by a packet at a link, i.e. $\theta_l(c_l, y_l)$, $l \in L$, has been predominantly based on approximation of links as independent $M/D/1$ queues. This approach stems from the Kleinrock independence approximation [5], which is in principle based on assumptions that the traffic arrives at network entry points according to a Poisson process, and the network is densely connected. Using this approach the average

packet delay $\theta_l(c_l, y_l)$ can be estimated as

$$\begin{aligned}\theta_l(c_l, y_l) &= \frac{1}{c_l} + \frac{y_l}{2c_l(c_l - y_l)}, \quad \forall l \in L \\ &\approx \frac{1}{c_l} + \frac{1}{2(c_l - y_l)}, \quad \forall l \in L\end{aligned}\tag{2.6}$$

The approximation (2.6) is based on assumption that at optimality y_l is close to c_l for all $l \in L$. Given approximation (2.6), the optimisation problem (2.1)-(2.5) is convex, and can be solved, for example, using primal decomposition [8] as follows

$$\max_{\mathbf{c}} \tilde{U}(\mathbf{c}) \quad \text{subject to} \quad (2.3) \tag{2.7}$$

where

$$\tilde{U}(\mathbf{c}) = \max_{\mathbf{x}} \sum_{s \in S} f_s(x_s) \quad \text{subject to} \quad (2.2), (2.4) \text{ and } (2.5) \tag{2.8}$$

Subproblem (2.8) can be solved using primal or dual algorithms proposed in [24]. Moreover, by standard convex programming results, \tilde{U} is concave (Proposition 3.4.3 in [4]), and therefore the set of optimal Lagrange multipliers associated with constraints (2.2) and (2.4) in subproblem (2.8) is the subdifferential of \tilde{U} (Section 5.4.4 in [4]). This property can then be used to develop algorithms for solving (2.7), when duality-based approaches are used to solve (2.8).

However, the $M/D/1$ queue approximation of links has several flaws. Firstly, the key assumptions behind the $M/D/1$ approximation (2.6) do not hold since the traffic at entry points are regulated by the rate controller and are deterministic, multi-hop networks are composed of mostly disjoint paths which comprise serial links, and the traffic entering the links can be further regulated to limit its burstiness [10, 11]. The delay caused by the burstiness of the arriving traffic at each link can therefore be assumed to be negligible and consequently the average delay a packet experiences at equilibrium is primarily a function of number of packets in the system at equilibrium, which is determined by the dynamics of the rate control and scheduling algorithms at their transient state. Secondly, in the approximation (2.6), as traffic rates at links approach their capacities, their delays grow exponentially. This implies that at optimality links are not efficiently utilised, in order to ensure bounded delay.

2.5 Conclusions

In this chapter the problem is formulated as a network utility maximisation problem. Furthermore, it is shown that when delay constraints are modelled based on approximation of links as $M/D/1$ queues, the problem is convex and can be solved by appropriate decomposition, and using the previously proposed algorithms in the literature.

However, the key assumption behind $M/D/1$ queue delay model, i.e. Poisson arrival of packets at links, does not hold in the problem under consideration, as delay is mainly determined by the transient behaviour of the rate control and scheduling algorithms. As such, the estimated delay in such models is inaccurate. Moreover, such delay models lead to inefficient utilisation of links at optimality since their estimated delay grows exponentially as link flow rates approach their capacities.

In next chapter, it will be shown that the predominant network utility maximisation approaches to support delay sensitive traffic, also have similar limitations.

Chapter 3

Related Work

3.1 Introduction

In this chapter the predominant approaches that aim to address the requirements of delay-sensitive traffic using the NUM framework, as well as their limitations, are described. In Section 3.2, the two approaches to model delay-sensitive traffic, namely, representation as non-concave utility functions and hard constraints specification, are explained. In Section 3.3, the well-known distributed solutions for the problem of joint rate control and scheduling for elastic traffic, are presented. Various approaches that tackle the same problem but for inelastic traffic, including minimising delay using virtual data rates, minimising network congestion, and minimising network distortion for video traffic, are introduced in Section 3.4. Finally, distributed rate control algorithms for networks with heterogeneous traffic, for cases where delay-sensitivity is modelled as a concave, as well as a non-concave function, are discussed in Section 3.5.

3.2 Modelling Delay-Sensitive Traffic

3.2.1 Representation as Non-Concave Utility Functions

In [34] it is argued that network performance should be evaluated on the basis of the degree to which the network satisfies the service requirements of user's applications, rather than in terms of network-centric measures such as link utilisation, dropped packets and so on. Let the vector s_i describe the service provided to the i th application or user, which contains all relevant measures like delay and throughput. The notion of utility function U_i is then defined as the mapping from the vector s_i into the performance of the application. The utility function describes how the performance of an application depends on the delivered service. The network design goal is subsequently defined as to maximise the performance of all applications, or in other words, to maximise the the sum of the utilities, also referred to as *efficacy*.

Based on the simplified assumption that the service can be merely described in terms of bandwidth, the shape of the utility functions for common classes of applications are described (Figure 3.1). Traditional data applications which are tolerant of delay have a diminishing marginal rate of performance enhancement as bandwidth is increased. So the utility function of such applications is strictly concave everywhere. The network efficacy is always maximised by admitting all users. Such applications are referred to as elastic applications. Hard real-time applications, on the other hand, require that data packets data arrive within a specified delay bound and perform very badly if the packets arrive later than this bound. The performance remains constant for bandwidths beyond the critical level needed to meet the required delay bound, but falls shapely for bandwidths below the critical level. The utility of these application looks like a step function. A network with such applications requires admission control to ensure the required bandwidth for these applications.

Delay-adaptive real-time applications like most current audio and video applications are rather tolerant of occasional delay-bound violations and packet loss. However, they still have an intrinsic bandwidth requirement as they cannot adapt

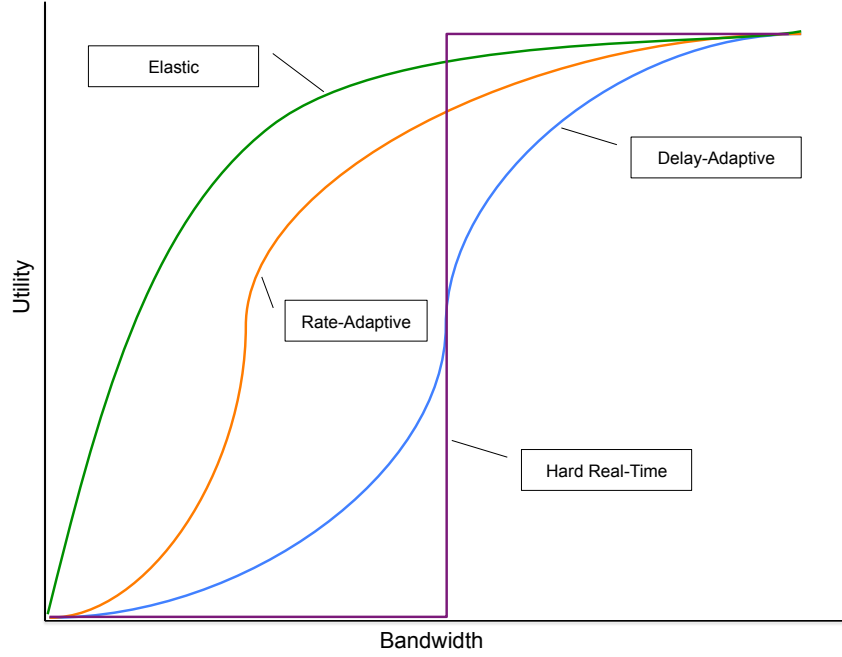


Figure 3.1: Utility functions for common classes of applications [34]

their data generation rate to the network congestion. As such, their performance degrades severely as soon as the bandwidth drops below the intrinsic rate, although not as sharply as with the hard real-time applications. The shape of the utility function is very similar to the utility function of the hard real-time applications, specially it is convex but not concave in the neighbourhood around zero. This implies that network can be overloaded with these applications at some point and hence has to use admission control to maximise its efficacy.

Rate-adaptive real-time applications adjust their transmission rate according to network congestion. Thus, their performance depends completely on the signal quality. At high bandwidths the marginal increase in utility as a result of additional bandwidth is very slight since the signal quality is much better than human need. Similarly, at very low bandwidths the marginal increase in utility as a result of additional bandwidth is very slight as the signal quality is very low. Similar to the utility functions of the delay-adaptive applications, these utility functions are

convex but not concave in the neighbourhood around zero and so the network can be overloaded with these applications.

It is finally concluded that for a network supporting only elastic traffic, the efficacy is maximised by admitting all flows. However, when there are also real-time traffic then the efficacy is maximised by rejecting some flows.

Rate Control for Elastic Traffic

For a network with only elastic traffic, the rate allocation problem, i.e. the NUM problem (2.1), (2.2), and (2.5) with fixed link data rates \mathbf{c} , becomes a convex optimisation problem, since all the utility functions are strictly concave. Thus its optimal solutions are global and the duality gap is zero [6]. This leads to the canonical distributed price-based rate control algorithm through solving the dual problem

$$\min_{\boldsymbol{\lambda} \geq 0} D(\boldsymbol{\lambda}) \quad (3.1)$$

where

$$\begin{aligned} D(\boldsymbol{\lambda}) &= \max_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\lambda}) \quad \text{subject to} \quad (2.5) \\ &= \max_{\mathbf{x}} \sum_{s \in S} f_s(x_s) - \boldsymbol{\lambda}^T (R\mathbf{x} - \mathbf{c}) \quad \text{subject to} \quad (2.5) \end{aligned} \quad (3.2)$$

and $\boldsymbol{\lambda}$ is the vector of Lagrange multipliers (link prices) associated with constraint (2.2). The dual function (3.2) can be decomposed into individual source problems as follows

$$\begin{aligned} D(\boldsymbol{\lambda}) &= \sum_{s \in S} \left(\max_{\mathbf{x}^s \geq 0} f_s(x_s) - \boldsymbol{\lambda}^T R^s \mathbf{x}^s \right) + \boldsymbol{\lambda}^T \mathbf{c} \\ &= \sum_{s \in S} D^s(\boldsymbol{\lambda}) + \boldsymbol{\lambda}^T \mathbf{c} \end{aligned} \quad (3.3)$$

The dual problem (3.1) is convex [6], and can be solved using the following sub-gradient method [30]

$$\lambda_l(t+1) = [\lambda_l(t) + \beta (R_l \mathbf{x}(\boldsymbol{\lambda}(t)) - c_l)]^+ \quad \forall l \in L \quad (3.4)$$

or in continuous-time form [35]

$$\dot{\lambda}_l = \beta [R_l \mathbf{x}(\boldsymbol{\lambda}) - c_l]_{\lambda_l}^+ \quad \forall l \in L \quad (3.5)$$

where $\beta > 0$, R_l is the l th row of R , and $\mathbf{x}(\boldsymbol{\lambda}(t))$ is the solution of (3.3) given $\boldsymbol{\lambda}(t)$, i.e.

$$\mathbf{x}^s(\boldsymbol{\lambda}) = \arg \max \left(f_s(x_s) - \boldsymbol{\lambda}^T R^s \mathbf{x}^s \right) \quad \forall s \in S \quad (3.6)$$

Furthermore, $[\cdot]^+$ denotes the projection on the set $\Lambda = \{\boldsymbol{\lambda} | \boldsymbol{\lambda} \geq \mathbf{0}\}$, and $[g(x)]_x^+$ is defined by

$$[g(x)]_x^+ = \begin{cases} g(x) & x > 0 \\ \max(g(x), 0) & x = 0 \end{cases}$$

In (3.4) or (3.5), each link $l \in L$ updates its price λ_l in proportion to the difference between its aggregate flow rate $R_l \mathbf{x}(\boldsymbol{\lambda}(t))$ and its data rate c_l . In (3.6), each source $s \in S$ adjusts its path rates \mathbf{x}^s according to its current path prices. It is shown that with appropriate choice of β , $\boldsymbol{\lambda}(t)$ converges to the dual optimal solution $\boldsymbol{\lambda}^*$ as $t \rightarrow \infty$. Given in the primal problem (2.1), (2.2), and (2.5), f_s , $s \in S$ are strictly concave, $\{x_s(\boldsymbol{\lambda}(t))\}$ converges to the primal optimal $\{x_s^*\}$. It then follows that in the case of single-path routing $\mathbf{x}(\boldsymbol{\lambda}(t))$ also converge to the primal optimal solution \mathbf{x}^* . The link price algorithm (3.4) or (3.5), and rate control algorithm (3.6) can be performed distributively by individual links and sources, respectively.

Rate Control for Inelastic Traffic

When the network supports a mixture of elastic and Inelastic traffic, the NUM problem (2.1), (2.2), and (2.5) with fixed link data rates \mathbf{c} , becomes a case of non-convex optimisation problem, since the utility functions for inelastic traffic are non-concave or non-smooth. Optimisation of the NUM problem with non-concave utility functions is generally difficult, as local optimum may not be a global optimum and the duality gap can be strictly positive. In this case the dual problem (3.1) is not equivalent to the primal problem anymore, and consequently the canonical distributed algorithm (3.4)-(3.6) may fail to converge to the primal optimal solution, or even a feasible rate allocation.

However, it is proved in [9] that for a non-concave NUM problem the canonical distributed algorithm (3.4)-(3.6) converges to a globally optimal rate allocation if the price-based rate allocation $\mathbf{x}^*(\boldsymbol{\lambda})$ is continuous at the optimal prices $\boldsymbol{\lambda}^*$.

Furthermore, continuity of $\mathbf{x}^*(\boldsymbol{\lambda})$ at at least one of the optimal prices $\boldsymbol{\lambda}^*$ is a necessary condition for the canonical distributed algorithm to converge to a globally optimal rate allocation. Although it is generally difficult to test for continuity of the price-based rate allocation for non-concave utility maximisation, for many types of non-concave utility functions (e.g. sigmoidal functions), it is easy to characterise the set of $\boldsymbol{\lambda}$ at which $\mathbf{x}^*(\boldsymbol{\lambda})$ is discontinuous. Thus, in [9] a method is also presented to bound the range of $\boldsymbol{\lambda}^*$ which can then be used to verify whether it excludes the points of discontinuity and hence to ensure the continuity of the price-based rate allocation.

The alternative approach proposed in [23] considers the particular case where services are classified into two types based on the shape of the utility function: traditional data services whose elasticity is modelled by a concave utility function, and delay and rate sensitive services (e.g. streaming video and audio) whose elasticity is modelled by a sigmoidal-like utility function. An increasing function $f(x)$ is called a *sigmoidal-like function*, if it has one inflection point x_0 , and $\frac{d^2f(x)}{dx^2} > 0$, for $x < x_0$, and $\frac{d^2f(x)}{dx^2} < 0$, for $x > x_0$. It is first shown that approximation of a sigmoidal-like utility function with a concave function and using algorithms developed for concave utility functions could result in a highly inefficient solution. Next, it is proved that when there are users with sigmoidal-like utility functions, the canonical distributed algorithm may cause link congestion without convergence. To avoid the congestion in the network some users have to be interrupted, and given that there is no central authority in the internet, a ‘self-regulating’ algorithm is proposed whereby each user ‘self-regulates’ its access to the network based on the local information. The proposed algorithm is then shown to converge to a rate allocation that does not lead to congestion when users with sigmoidal-like utility functions are present. Furthermore, the resulting rate allocation is asymptotically optimal, that is, it is a good approximation of global optimal rate allocation when there are many users in a system with large capacity and the number of users that stop transmitting data due to the ‘self-regulating’ property has vanishing proportion.

3.2.2 Representation as Hard Constraints

In [9] a different modelling approach for inelastic traffic is presented where instead of using non-concave utility functions, inelastic traffic are modelled explicitly as hard constraints and objective functions of the NUM problem. Delay-sensitive traffics are classified into three types. The R-type traffic indexed by r represents real-time applications, such as real-time IP, which require constant playback rate and a fixed requested playback starting time. The R-type traffic is specified as follows

- Playback rate is required to be a constant of v_r bits per time unit.
- For each source an admission decision a_r is made. If the flow is admitted a constant utility \bar{U}_r is gained.
- The optimisation problem for R-type traffic is an admission control problem with optimisation variables are $a_r \in \{0, 1\}$.

The B-type traffic indexed by b represents streaming applications which require constant playback rate, but have a flexible playback starting time. It is assumed that a playback buffer at the receiver can absorb fluctuations of the source rate to some extent. The B-type traffic is specified as follows

- Playback rate is required to be a constant of v_b bits per time unit.
- The optimisation variables are the actual playback start time w_b and the rate allocation over time $x_b(t)$. The optimisation problem for B-type traffic is joint problem of scheduling and rate allocation over time slots.
- To guarantee the constant playback, receiver buffers should not be depleted during playback, i.e.

$$\sum_{t=0}^{t_0} x_b(t) \geq (t_0 - w_b)v_b, \quad \forall t_0 > w_b$$

- The utility is $U_b(w_b)$ where U_b is non increasing, since users prefer to start the playback as early as possible.

The D-type traffic indexed by d represents general delay-sensitive traffic whose utility depends on the transient behaviour of rate allocation. The D-type traffic is specified as follows

- The optimisation variables are the rate allocation over time $x_d(t)$. The optimisation problem for D-type traffic is a rate allocation problem over all time slots.
- The utility is assumed to be $\sum_t U_{d,t}(x_d(t))$, where $U_{d,t}$ are concave and increasing.

A simplified version of NUM problem for rate allocation among the three types of inelastic flows on a single link is then given by

$$\begin{aligned}
& \max_{\mathbf{x}_D, \mathbf{a}_R, \mathbf{x}_B, \mathbf{w}_B} \quad \sum_{d \in D} \sum_{t=0}^T U_{d,t}(x_d(t)) + \sum_{r \in R} a_r \bar{U}_r + \sum_{b \in B} U_b(w_b) \\
& \text{subject to} \quad \sum_{d \in D} x_d(t) + \sum_{r \in R} a_r v_r + \sum_{b \in B} x_b(t) \leq c \quad \forall t \\
& \quad \sum_{t=0}^{t_0} x_b(t) \geq (t_0 - w_b) v_b \quad \forall t_0 > w_b, \quad \forall b \in B \\
& \quad x_d(t) \geq 0, x_b(t) \geq 0 \quad \forall d \in D, b \in B, \quad \forall t \\
& \quad a_r \in \{0, 1\} \quad \forall r \in R \\
& \quad w_b \in [0, T] \quad \forall b \in B
\end{aligned}$$

The first constraint ensures that the traffic flow on the link does not exceed its capacity. The second constraint ensures that receiver buffers are not depleted during playback. The above formulation applies to the single link case but can be readily generalised to arbitrary network topologies. Furthermore, to simplify the model, it is assumed that all requested flow starting times are time 0, each receiver playback buffer is infinitely large and all flows have infinite backlogs.

It is shown that the above optimisation problem can be decomposed into individual source problems and link problems which can then be solved using an algorithm similar to the canonical distributed algorithm for elastic traffic. However, the proposed algorithm has two major limitations. Firstly, congestion prices

have be generated using an iterative algorithm for every time slot along the temporal dimension (as well as for every link in the case of time-sensitive traffic). Secondly, optimal admission decision, playback time decision, and rate allocation cannot be made until the end of the entire period $t = 0, \dots, T$. For the special case where there are only inelastic real-time traffic and elastic TCP traffic present, the canonical distributed algorithm can be used for price update and elastic source rate control, but admission decision of real-time flows can only be made after the equilibrium price is reached. A price-based admission control heuristics is then proposed to avoid the delay associated with optimal admission decision.

3.2.3 Comparison of the Modelling Approaches

Representation of delay-sensitivity via the utility functions offers several advantages over the hard-constraint representation. Firstly, since most applications on the internet have some degree of elasticity to the allocated rate, utility functions capture more accurately the level of user satisfaction or QoS at the allocated rate. Secondly, the elasticity modelled by the utility function can then be exploited through rate control to maximise network efficacy, using a NUM framework. Thirdly, as discussed previously, for many types of non-concave utility functions, it is easy to verify that the canonical distributed algorithm (3.4)-(3.6) converges to the optimal rate allocation. In addition, when some users have sigmoidal-like utility functions, the canonical distributed algorithm combined with the distributed ‘self-regulating’ algorithm proposed in [23] converges to a rate allocation that is asymptotically optimal. On the other hand, the hard-constraint modelling of delay-sensitive traffic presented in Section 3.2.2 leads to more computationally complex solutions.

However, the characterisation of delay-sensitivity via (non-concave) utility functions as described in Section 3.2.1 is based on the simplified assumption that the delay can be merely described in terms of source’s allocated rate. While this is true for delay between consecutive packets arrival at the receiver, it is not a correct assumption for packets end-to-end queueing delay. As explained in Section 2.4, end-to-end queueing delay is dependent on the link congestion levels which are

determined by dynamics of the rates of flows passing through individual links and link data rates. Thus the utility functions in the form described in Section 3.2.1 can only capture the sensitivity of an application to the allocated data rate, but not to packet end-to-end delay. Similarly, end-to-end queueing delay is ignored in formulating the constraints and objective functions of delay-sensitive traffic presented in Section 3.2.2. The coming sections describe the predominant approaches that attempt to address end-to-end queueing delay requirements of delay-sensitive traffic, using the NUM framework.

3.3 Joint Rate Control and Scheduling for Elastic Traffic

The problem of joint rate control and scheduling for elastic traffic, i.e. the optimisation problem (2.1), (2.2), (2.3) and (2.5), has been extensively studied (see e.g. [29] and references therein and [7]). Dual optimisation-based approach has been the preferred solution strategy for this problem since it enables decomposition of the problem into the rate control and scheduling ‘layers’ coupled loosely through ‘link prices’. Specifically, the dual problem is given by

$$\min_{\boldsymbol{\lambda} \geq 0} D(\boldsymbol{\lambda}) \quad (3.7)$$

where

$$\begin{aligned} D(\boldsymbol{\lambda}) &= \max_{\mathbf{x}, \mathbf{c}} L(\mathbf{x}, \mathbf{c}, \boldsymbol{\lambda}) \quad \text{subject to} \quad (2.3) \text{ and } (2.5) \\ &= \max_{\mathbf{x}, \mathbf{c}} \sum_{s \in S} f_s(x_s) - \boldsymbol{\lambda}^T (R\mathbf{x} - \mathbf{c}) \quad \text{subject to} \quad (2.3) \text{ and } (2.5) \end{aligned} \quad (3.8)$$

and $\boldsymbol{\lambda}$ is the vector of Lagrange multipliers associated with constraint (2.2). Using the shadow price interpretation of Lagrange variables [6], $\boldsymbol{\lambda}$ can also be interpreted as the link data rate prices. The optimisation problem (3.8) can be decomposed into the following rate control and scheduling subproblems, respectively

$$D_1^s(\boldsymbol{\lambda}) = \max_{\mathbf{x}^s \geq 0} f_s(x_s) - \boldsymbol{\lambda}^T R^s \mathbf{x}^s \quad \forall s \in S \quad (3.9)$$

and

$$D_2(\boldsymbol{\lambda}) = \max_{\mathbf{c} \in u(\mathbf{p}), \mathbf{p} \in \Pi} \boldsymbol{\lambda}^T \mathbf{c} \quad (3.10)$$

The dual problem (3.7) can be solved using the subgradient method as follows

$$\lambda_l(t+1) = [\lambda_l(t) + \beta (R_l \mathbf{x}(\boldsymbol{\lambda}(t)) - c_l(\boldsymbol{\lambda}(t)))]^+ \quad \forall l \in L \quad (3.11)$$

where $\beta > 0$, R_l is the l th row of R , and $\mathbf{x}(\boldsymbol{\lambda}(t))$ and $\mathbf{c}(\boldsymbol{\lambda}(t))$ are the solutions of (3.9) and (3.10) given $\boldsymbol{\lambda}(t)$, respectively. Here $[\cdot]^+$ denotes the projection on the set $\Lambda = \{\boldsymbol{\lambda} | \boldsymbol{\lambda} \geq \mathbf{0}\}$. In (3.11), each link $l \in L$ updates its price λ_l in proportion to the difference between its aggregate flow rate $R_l \mathbf{x}(\boldsymbol{\lambda}(t))$ and its data rate $c_l(\boldsymbol{\lambda}(t))$. The rate control subproblem (3.9) and the scheduling subproblems (3.10) are coupled via link prices $\boldsymbol{\lambda}(t)$. In (3.9), each source $s \in S$ adjusts its path rates \mathbf{x}^s according to its path prices. In (3.10), link data rates \mathbf{c} are updated based on the link prices. The link price algorithm (3.11) and rate control algorithm (3.9) can be performed in a distributed fashion by individual links and sources, respectively.

3.3.1 Scheduling Solution Approaches

The scheduling problem (3.10) is a computationally complex problem in general, since $u(\mathbf{p})$ is not concave in many cases and as a result convex programming methods cannot be used. Given the fact that link prices $\boldsymbol{\lambda}(t)$ are updated at every timeslot and therefore (3.10) has to be solved at every timeslot, finding an efficient, simple and distributed solution becomes extremely vital. Cases where the scheduling problem (3.10) is solvable include [29]:

- *node-exclusive interference model*, where each wireless node can only communicate with one other node at any time. This model represents Bluetooth-like networks with high accuracy and is a reasonable approximation to frequency-hopping code-division multiple-access (FH-CDMA) systems. In this case the scheduling problem corresponds to a maximum-weighted-matching problem, which has polynomial-time complexity.

- *low-SINR model*, where the data rate of each link is a linear function of its SINR. This model approximates CDMA systems with a moderate processing gain. In this case, optimal scheduling is shown to have the property where each node either transmit at maximum power to only one other node, or does not transmit at all. This property substantially reduces the search space for optimal scheduling but the problem still has exponential complexity in the number of nodes.

In certain cases including the *high-SINR* model, where the data rate of each link is a logarithmic function of its SINR; the *low-SINR* model, and the single channel Aloha networks, the function $u(\mathbf{p})$ can be transformed into a concave function after some change of variables [29]. Standard convex programming methods can then be used to solve the transformed problem.

An alternative approach to compute the exact solution for the scheduling problem (3.10) is to instead find suboptimal solutions that are simpler and enable distributed implementation[27, 29]. In [27] the fairness and efficiency of a cross-layer solution using a class of imperfect scheduling policies referred to as S_γ -policies are reviewed. Roughly speaking, an S_γ -policy can guarantee a minimum capacity region of $\gamma\Lambda$, where $\gamma \in (0, 1]$, and the capacity region Λ is the largest set of transmission rates \mathbf{x} such that for every $\mathbf{x} \in \Lambda$ there exists some scheduling policy that can stabilise the network. Although for the case where the user population is fixed only a weak fairness property can be shown, the stability region is shown to be at least $\gamma\Lambda$ when user population varies according to a stochastic process. The class of maximal scheduling policies, which are very simple and can be implemented in a distributed fashion, are then shown to be S_γ -policies with $\gamma = \frac{1}{2}$ under the node-exclusive interference model.

3.4 Joint Rate Control and Scheduling for Delay-Sensitive Traffic

3.4.1 Minimising Delay Using Virtual Data Rates

Since the algorithm (3.11) couples the link prices to their average queue lengths, it may lead to large queue lengths and hence large queueing delays at the equilibrium. As suggested in [26, 32, 22], this can be avoided by using the slightly smaller ‘virtual’ link data rates in (3.11) instead of the actual link data rates. Specifically, $c_l(\boldsymbol{\lambda}(t))$, $l \in L$ in (3.11) is replaced by $\rho c_l(\boldsymbol{\lambda}(t))$, where ρ is a positive factor slightly smaller than 1. While the modified algorithm still leads to the link prices close to their optimal level, it results in zero equilibrium queue lengths, since links traffic loads are slightly less than their actual data rates at equilibrium. Main disadvantages of this approach are that it does not completely utilise network capacity and provides no control over the transient behaviour of packet delays.

3.4.2 Minimising Network Congestion

In [40] the problem of joint optimisation of link capacities (data rates) and flow assignment for delay-sensitive applications with focus on live video streaming is considered. It is assumed that the data rate of each link $l \in L$ is given by

$$c_l = W \log_2 \left(1 + \frac{SINR_{\mathbf{p},l}}{\Gamma} \right) \quad (3.12)$$

where W is the system bandwidth, Γ is a constant determined by the BER requirement and the coding scheme, and $SINR_{\mathbf{p},l}$ is SINR at the receiving node of link l given the power assignment vector \mathbf{p} . To reduce the computational complexity, the transmission power of the transmitting nodes are fixed at their maximum level. Furthermore, if a particular vector of link data rates \mathbf{c} can be generated by appropriate time division of the other link data rates, it is removed. Links with gains below a certain threshold are also removed to prevent long range communications and force (upper-layer) routing algorithms to use multi-hop routing to reach the destination.

The primary objective is to find a resource allocation strategy that supports maximum data rates and yields minimum end-to-end delay; however, for general queueing systems, this leads to an intractable problem formulation. Hence, an alternative problem formulation is proposed where the network congestion, as a measure of delay experienced by packets, is minimised while allowing communication between source and destinations at a given data rate. Specifically, the network congestion is defined as the maximum link utilisation over all links

$$\Delta(\mathbf{c}, \mathbf{y}) = \max_{l \in L} \frac{y_l}{c_l} \quad (3.13)$$

The alternative optimisation problem is then given by

$$\begin{aligned} \min_{\mathbf{c}, \mathbf{y}} \quad & \Delta(\mathbf{c}, \mathbf{y}) \\ \text{subject to} \quad & (2.2), (2.3), (2.5) \\ & \mathbf{y} = R\mathbf{x} \\ & x_s \geq m_s \quad \forall s \in S \end{aligned} \quad (3.14)$$

The objective function (3.13) is a quasi-convex function of \mathbf{c} and \mathbf{x} . The above optimisation problem can then be solved, for example, by a bisection algorithm that involves solving a sequence of convex feasibility problems [6]. The optimal flow assignment however may not be very practical, since it does not directly indicate the set of paths between sources and destinations, and it may use a large number of links. Therefore, in a recursive process, the k paths carrying the most traffic are selected from the optimal solution, and then the optimisation problem is resolved by constraining the flows to the selected paths. Experimental results indicate that the proposed cross-layer optimisation approach results in significant improvement in supported data rate and video quality, compared with a method based on oblivious layers. Evidently, this approach does not aim to efficiently utilise links.

3.4.3 Minimising Total Distortion for Video Transmissions

In [41] the problem of optimal rate allocation for multi-stream video transmission over wireless ad hoc networks is considered. The main focus is on designing a

distributed rate allocation algorithm that minimises total distortion of all video streams and can easily adopt to fluctuations in the wireless network conditions. The following convex optimisation problem is then defined

$$\begin{aligned} \min_{\mathbf{x} \geq 0} \quad & \sum_{s \in S} w_s D_s(x_s) \\ \text{subject to} \quad & R\mathbf{x} \leq \alpha \mathbf{c} \end{aligned} \tag{3.15}$$

where w_s is the relative importance of source $s \in S$, and $D_s(x_s)$ is the video distortion for source $s \in S$ defined by

$$D_s(x_s) = D_{0,s} + \frac{\theta_s}{x_s - x_{0,s}}$$

where the parameters $D_{0,s}$, θ_s and $x_{0,s}$ are estimated from trial encodings. The constraint in (3.15) requires that the traffic rate entering each link to be below the actual link data rate by a margin determined by the scaling factor $\alpha < 1$. It is assumed that link data rates \mathbf{c} are determined by a media access control mechanism, which is dependent on the traffic pattern on each link. Optimal rates are then computed using a subgradient algorithm where link prices updated similarly to (3.11), in which at every step the actual data rate (capacity) and flow rates at each link are estimated from the observed packet arrival and departure times, averaged over many packets. It is shown using simulation that despite the inaccuracies associated with the proposed scheme, the achieved source rates are similar to those obtained from exhaustive search.

3.4.4 Providing Bounded Delay for Traffic with Elastic Bandwidth Requirements

This thesis extends the ideas presented in author's previous papers [19, 20]. In both papers the correlation between optimal link prices and equilibrium link average queueing delays in duality-based rate control algorithm is exploited in order to provide bounded average end-to-end queueing delay. In [19], an approach is presented first in which lower bounds on sources' transmission rates are derived in order to ensure the required bounded delay. This approach inevitably entails

admission control. In the second approach, an alternative formulation is introduced where the delay constraint is omitted and instead the utility function for each source is multiplied by a weight factor. The proposed solution comprises a scheduling algorithm incorporating a duality-based rate control algorithm at its inner layer, and an algorithm that dynamically adjusts sources' weights to ensure the required bounded delay. In this thesis the latter approach is developed by designing a new scheduling algorithm and delay regulator that regulate average queueing delay with high accuracy and performance. Moreover, a complete analysis of the stability of the proposed algorithms is provided. In [20] the sensitivity of optimal path prices for each source to the variation of its weight factor is analysed, and a delay regulator is presented that is integrated into the duality-based rate control and scheduling algorithms given in [29]. Here the sensitivity analysis results in [20] is used to develop a solution that regulates average queueing delay with higher accuracy and performance.

3.5 Rate Control for Heterogeneous Traffic

3.5.1 Maximising Utility as a Function of Rate and Delay - The Concave Case

In [24] the congestion control problem in networks supporting traffic with various levels of rate, delay and packet loss sensitivity, is studied. The proposed approach is based on incorporating the requirements for rate, delay and packet loss in the utility function of sources. It is assumed that each source transmits only one flow using a fixed path, and that link data rates \mathbf{c} are fixed. The average delay experienced by a packet at link l is given by $\theta_l(y_l)$, where θ_l are assumed to be positive, increasing and convex for all $l \in L$. The average packet delay for each source $s \in S$ is then given by $\sum_{l \in L} R_l^s \theta_l(y_l)$. The utility of each source $s \in S$ is subsequently defined by

$$U_s = a_s f_s(x_s) - b_s \sum_{l \in L} R_l^s \theta_l(y_l)$$

where coefficients a_s and b_s indicate the degree of sensitivity of the traffic to rate and delay, respectively. In particular, $(a_s = 0, b_s > 0)$ for delay-sensitive traffic with fixed rate requirement like VoIP, $(a_s > 0, b_s > 0)$ for delay-sensitive and rate-sensitive traffic like real-time data, and $(a_s > 0, b_s = 0)$ for rate-sensitive traffic with no delay requirement like file-downloading. The optimisation problem is

$$\begin{aligned}
 & \max_{\mathbf{x}, \mathbf{y}} \quad \sum_{s \in S} U_s(x_s, \mathbf{y}) \\
 & \text{subject to} \quad R\mathbf{x} = \mathbf{y} \\
 & \quad \mathbf{y} < \mathbf{c} \\
 & \quad \mathbf{x} \geq 0
 \end{aligned} \tag{3.16}$$

It is shown that similar to the basic congestion control problem for elastic traffic, the above alternative optimisation problem can be solved using either primal algorithms [35] with link prices $p_l(y_l(t)) = (\sum_{s \in S} R_l^s b_s) \theta'_l(y_l(t))$, or dual algorithms [35] where $c_l(\lambda_l(t)) = \theta_l'^{-1}\left(\frac{\lambda_l(t)}{\sum_{s \in S} R_l^s b_s}\right)$.

The analysis is then applied to networks with mixed voice and data traffic, for the cases when priority queueing is used and when it is not. When priority queueing is used, voice packets are given higher priority than data packets. The arrival processes of voice and data are assumed to be independent, Poisson, and independent of the service times. Two separate queues are maintained for voice packets and data packets, respectively. It is assumed that each link is independent and that the arrival processes of voice and data at each link are also independent. Let subscripts D and V denote data and voice, respectively. In this case the average delay of a voice packet for source s^V is

$$\delta_{s^V}(\mathbf{y}^D) = \sum_{l \in L} R_l^{s^D} \left(\frac{K}{c_l} + \frac{K}{2c_l} \frac{y_l^V + y_l^D}{c_l - y_l^V} \right)$$

Moreover, the average delay of a data packet for source s^D is

$$\delta_{s^D}(\mathbf{y}^D) = \sum_{l \in L} R_l^{s^D} \left(\frac{K}{c_l} + \frac{K}{2(c_l - y_l^V)} \frac{y_l^V + y_l^D}{c_l - y_l^V - y_l^D} \right)$$

where \mathbf{y}^V is fixed for voice traffic. The utility function for the source of voice traffic s^V is set as a function of its R-factor denoted by $R_{s^V}^{fac}$. $R_{s^V}^{fac}$ is typically a convex

function of average packet delay δ_{s^V} and packet loss probability given by $\psi_{s^V} = p_{s^V} + \phi_{s^V}$, where p_{s^V} is the packet loss probability due to an unreliable wireless link, and ϕ_{s^V} is the probability that the delay of a packet exceeds its deadline d_s . Considering only packet loss as a result of unreliable links, i.e. $\psi_{s^V} = p_{s^V}$, $R_{s^V}^{fac}$ is then a linear function of average packet delay δ_{s^V} . The utility function for the source of data traffic s^D is set as the weighted sum of utility on throughput and utility on delay. The optimisation problem is formulated as

$$\begin{aligned}
& \max_{\mathbf{x} \geq 0} \quad \frac{v}{S^V} \sum_{s^V \in S^V} U_{s^V}(R_{s^V}^{fac}) + \frac{1-v}{S^D} \sum_{s^D \in S^D} U_{s^D}(x_{s^D} \rho_{s^D}, \delta_{s^D}) \\
& \text{subject to} \quad \mathbf{y}^V + \mathbf{y}^D \leq \mathbf{c} \\
& \quad \bar{R}_{s^V}^{fac} \leq R_{s^V}^{fac}(\delta_{s^V}) \quad \forall s^V \in S^V \\
& \quad \delta_{s^V}(\mathbf{y}^D) \leq \bar{\delta}_{s^V} \quad \forall s^V \in S^V \\
& \quad \delta_{s^D}(\mathbf{y}^D) \leq \bar{\delta}_{s^D} \quad \forall s^D \in S^D
\end{aligned} \tag{3.17}$$

where $v \in [0, 1]$, \bar{R}_{s^V} , $\bar{\delta}_{s^V}$ and $\bar{\delta}_{s^D}$ are constants, and $\rho_{s^D} = 1 - p_{s^D}$, where p_{s^D} is the end-to-end packet loss probability for source $s^D \in S^D$. The second constraint states that the R-factor of each voice traffic should not be less than the requested R-factor. The third and forth constraints impose upper bounds on the average end-to-end delay of voice and data packets, respectively. Based on the dynamics of the proposed dual algorithm, a distributed algorithm for solving the above optimisation problem is then proposed.

For the case with no priority queueing it is assumed that incoming packets on a link are stored in a queue and transmitted on FIFO basis. Furthermore, it is assumed that packet arrivals at entry points Poisson processes and each link can be approximated as an independent $M/D/1$ queue. In this case the average packet delay for source s is

$$\delta_s(\mathbf{y}^D) = \sum_{l \in L} R_l^s \left(\frac{K}{2c_l} + \frac{K}{2(c_l - y_l^V - y_l^D)} \right)$$

The optimisation problem and its distributed solution is subsequently derived similarly to the case with priority queueing.

In the general case where packet loss probability is given by $\psi_{s^V} = p_{s^V} + \phi_{s^V}$, it is hard to derive an exact formula for ϕ_{s^V} . By Markov inequality, $\frac{\delta_{s^V}}{d_s}$ provides

an upper bound for ϕ_{sv} . Using this upper bound instead of ϕ_{sv} in R_{sv}^{fac} , although the optimisation problem (3.17) is not convex in general, it is shown in [24] to be convex in the case of no priority queueing. Finding a distributed solution for this case however is still difficult. Numerical results in [24] show that the priority queueing improves both the R-factor of voice traffic and the throughput of data traffic, at the expense of the packet delay of data traffic.

However, the estimation of packet delay is based on approximation of links as independent $M/D/1$ queues which, as discussed in Section 2.3, stems from unrealistic assumptions and results in under-utilised links under optimal resource allocation.

3.5.2 Maximising Utility as a Function of Rate and Delay - The Non-Concave Case

In [36] a variant of the basic congestion control problem for elastic traffic [35] is studied where traffic sources are explicitly sensitive to delay as well as flows. Sources are heterogeneous with respect to their levels of sensitivity to both rate and delay. It is assumed that source $s \in S$ incurs a delay cost $h_s d$ per unit of flow rate, where d is the average end-to-end delay experienced by a packet. The utility of each source $s \in S$ is subsequently defined by

$$U_s = f_s(x_s) - h_s x_s \sum_{l \in L} R_l^s \theta_l(y_l) \quad (3.18)$$

When sources are homogeneous in their delay sensitivities, i.e. $h_s = h$, the optimisation problem (3.16) with U_s defined as in (3.18) becomes similar to the optimisation problem (3.1) in [35], given the price function $f_l(y_l)$ in equation (3.1) in [35] is replaced by C_l' , where $C_l(y_l) = h y_l \theta_l(y_l)$. In this case the objective function is strictly concave and hence the optimisation problem (3.16) has a unique optimal solution.

However, when sources are heterogeneous with respect to their delay sensitivities, the optimisation problem (3.16) with U_s defined as in (3.18) is shown to be non-concave in general and consequently may have several stationary points. A

dynamic rate control algorithm similar to the primal algorithm (3.8) in [35] is then considered, in which each source adjusts its rate in proportion to the difference between its marginal utility and its path price. The price of each link is set dynamically as the external effect of the flow that go through it at each time. It is shown that the considered rate control algorithm converges to a local maximum, which is a Nash equilibrium for the sources when charged the appropriate price, but never to a saddle point. A Nash equilibrium is a strategy profile (source transmission rates) \mathbf{x}^* such that no player (source) s can profit by unilaterally deviating from its strategy x_s^* , assuming every other player (source) \hat{s} follows its strategy $x_{\hat{s}}^*$ [31].

Two variants of the dynamic algorithm, one using fixed pricing and one using dynamic pricing based on total load at the link, are also shown to converge but generally not to the socially optimal points. It is concluded that dynamic rate control algorithms such as TCP may not be able to attain efficient rate allocations and levels of delay that are acceptable to diverse classes of traffic, in the absence of differentiated services. In this thesis it is assumed that the average packet end-to-end delay on each path is upper bounded in which case, as explained in Section 2.3, the problem can be formulated as convex optimisation problem (2.1)-(2.5), given approximation (2.6).

3.6 Conclusions

In this chapter, first, two prominent approaches to model the requirements of delay-sensitive traffic are introduced. In the first approach, it is assumed that delay can be described in terms of a traffic source's transmission rate, and QoS of a delay-sensitive traffic is typically represented as a non-concave utility function of its transmission rate. Since most current applications have some degree of elasticity to the allocated rate, this approach captures more accurately the QoS as a function of the allocated rate. Using the NUM framework, the modelled traffic elasticity can then be exploited to maximise network efficacy. Furthermore, for many types of non-concave utility functions, it is easy to verify that the canonical distributed algorithm, which solves the NUM problem for the elastic traffic, con-

verges to the optimal rate allocation. In the second approach, inelastic traffic are modelled explicitly as hard constraints and objective functions of the NUM problem, which normally leads to a computationally complex problem. Nevertheless, both approaches in their current form only characterise QoS in terms of packet inter arrival delay rather than end-to-end queueing delay.

Next, the prominent solution approach for the joint rate control and scheduling NUM problem for the elastic traffic is presented, which elegantly decomposes the problem into rate control and scheduling problems. The rate control problem is simply solved using the canonical distributed algorithm. The structure of the scheduling problem has also enabled the development of simple, efficient and distributed scheduling algorithms for many cases.

The predominant NUM approaches that aim to address end-to-end queueing delay requirements of delay-sensitive traffic have been mainly based on either reducing link utilisation, or approximation of links as $M/D/1$ queues. The former approach which includes using virtual data rates [26, 32] and minimising network congestion [40], normally leads to nearly zero queue lengths in the long term due to reduced link utilisation, but provides no control over the transient behaviour of packet delays. The latter approach which is adopted in [24] and [36] is based on assumptions that contrast with realistic scenarios. Moreover, it also results in under-utilised links under optimal resource allocation. These limitations motivate the research objectives stated in Section 1.3.4.

The elegance and simplicity of the solution algorithms, as well as the efficiency of the optimal solutions of the NUM problem when traffic QoS is characterised as a concave utility function of its transmission rate is the main motivation behind the proposed approach in this thesis. As will be described in the next chapters, the proposed approach is based on representation of delay-sensitive traffic QoS as a concave utility function of its transmission rate, whose shape is adjusted as the algorithm converges.

Chapter 4

Alternative Problem Formulation

4.1 Introduction

Given the limitations of the optimisation problem (2.1)-(2.5) discussed in Chapter 2, in this chapter an alternative formulation and its proposed solution is presented, on which the proposed approach for providing bounded delay will be based. In Section 4.2 the proposed alternative optimisation problem, in which the delay constraints are omitted and utility functions are multiplied by weight coefficients, is presented and its properties are examined. In Section 4.3 the scheduling representation of the alternative problem is introduced, which is based on vertical decomposition of the problem into master scheduling problem and the well-known multipath rate control with fixed link data rates subproblem. After deriving conditions for the uniqueness of the solutions of the multipath rate control subproblem, the proposed distributed algorithm for solving the scheduling problem is presented.

4.2 The Alternative Optimisation Problem

The proposed alternative optimisation problem is given by

$$\max_{x,c} \sum_{s \in S} w_s f_s(x_s) \quad \text{subject to} \quad (2.2), (2.3), \text{ and } (2.5) \quad (4.1)$$

where $w_s f_s(x_s)$ represents source s 's utility, or preference over transmission rate x_s . Compared with original problem (2.1)-(2.5), in the alternative problem (4.1), delay constraint (2.4) has been omitted and instead the utility function for each source $s \in S$ is multiplied by the weight parameter w_s . A geometric interpretation of w_s is that higher (respectively, lower) values of w_s results in higher (respectively, lower) marginal increase in source s 's preference or utility at a particular rate.

The dual problem for (4.1) is given by

$$\min_{\lambda \geq 0, \mu \geq 0} D(\boldsymbol{\lambda}, \boldsymbol{\mu}) \quad (4.2)$$

where

$$\begin{aligned} D(\boldsymbol{\lambda}, \boldsymbol{\mu}) &= \max_{\mathbf{x}, \mathbf{c}} L(\mathbf{x}, \mathbf{c}, \boldsymbol{\lambda}, \boldsymbol{\mu}) \quad \text{subject to} \quad (2.3) \\ &= \max_{\mathbf{x}, \mathbf{c}} \sum_{s \in S} w_s f_s(x_s) - \boldsymbol{\lambda}^T (R\mathbf{x} - \mathbf{c}) + \boldsymbol{\mu}^T \mathbf{x} \quad \text{subject to} \quad (2.3) \end{aligned} \quad (4.3)$$

and $\boldsymbol{\lambda}$ and $\boldsymbol{\mu}$ are the Lagrange multipliers associated with constraints (2.2) and (2.5), respectively. Using the shadow price interpretation of Lagrange variables [6], $\boldsymbol{\lambda}$ can also be interpreted as the link data rate prices.

Since optimisation problem (4.1) is convex and constraints (2.2) and (2.5) are affine, by Slater's theorem [6] the optimal duality gap is zero. Let $(\mathbf{x}^*, \mathbf{c}^*)$ and $(\boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$ be the primal and dual optimal solutions, respectively. Let also $\mathbf{q}^* = R^T \boldsymbol{\lambda}^*$. It then follows from Karush-Kuhn-Tucker (KKT) optimality conditions [6] that

$$w_s f'_s(x_s^*) - q_i^{s*} + \mu_i^{s*} = 0 \quad \forall i \in I_s, \quad \forall s \in S \quad (4.4)$$

$$\lambda_l^* (R_l \mathbf{x}^* - \mathbf{c}_l^*) = 0 \quad \forall l \in L \quad (4.5)$$

$$\mu_i^{s*} x_i^{s*} = 0 \quad \forall i \in I_s, \quad \forall s \in S \quad (4.6)$$

Equation (4.6) implies that $\mu_i^{s*} = 0$ for any $i \in I_s$ for which $x_i^{s*} > 0$. It then follows from (4.4) that

$$\begin{aligned} q_i^{s*} &= w_s f'_s(x_s^*) \quad i \in I_s, \quad x_i^{s*} > 0, \quad \forall s \in S \\ &\triangleq q_s^* \end{aligned} \quad (4.7)$$

which means that for each source $s \in S$ the values of q_i^{s*} associated with paths with positive flows are minimum and hence equal. Since the objective function in (4.1) is strictly concave with respect to $\{x_s\}$, $\{x_s^*\}$ is unique and it follows from (4.7) that \mathbf{q}^* is also unique. However, (4.1) is not strictly concave in either \mathbf{x} or \mathbf{c} , and hence neither \mathbf{x}^* or \mathbf{c}^* may be unique. Furthermore, given that $\mathbf{q}^* = R^T \boldsymbol{\lambda}^*$, and R may have linearly dependent rows, $\boldsymbol{\lambda}^*$ may not be unique in general.

4.3 Representation as a Scheduling Problem

Optimisation problem (4.1) can be alternatively presented as the following equivalent form

$$\max_{\mathbf{c}} U_{\mathbf{w}}(\mathbf{c}) \quad \text{subject to} \quad (2.3) \quad (4.8)$$

where

$$U_{\mathbf{w}}(\mathbf{c}) = \max_{\mathbf{x}} \sum_{s \in S} w_s f_s(x_s) \quad \text{subject to} \quad (2.2) \text{ and } (2.5) \quad (4.9)$$

The key feature of the alternative form (4.8) is the decomposition of the problem into master scheduling problem (4.8), and the well-known rate control subproblem (4.9) with fixed link data rates. In addition, $U_{\mathbf{w}}$ is concave by Proposition 3.4.3 in [4], and therefore, as shown in Section 5.4.4 in [4], the set of optimal Lagrange multipliers associated with constraint (2.2) in subproblem (4.9) is the subdifferential of $U_{\mathbf{w}}$.

The dual problem for (4.9) is similar to (4.2) and (4.3) with fixed link data rates \mathbf{c} . Consequently, KKT conditions (4.4)-(4.6) and Equation (4.7) also hold for problem (4.9). Let $\mathbf{x}(\mathbf{c})$ and $(\boldsymbol{\lambda}(\mathbf{c}), \boldsymbol{\mu}(\mathbf{c}))$ be the primal and dual optimal solutions of (4.9) given \mathbf{c} , respectively. Let also $\mathbf{q}(\mathbf{c}) = R^T \boldsymbol{\lambda}(\mathbf{c})$. Since the objective function in (4.1) is strictly concave with respect to $\{x_s\}$, $\{x_s(\mathbf{c})\}$ is unique and it follows from (4.7) that $\mathbf{q}(\mathbf{c})$ is also unique. However, as in the case of problem (4.1), $\boldsymbol{\lambda}(\mathbf{c})$ is not unique in general. Hence $\boldsymbol{\lambda}(\mathbf{c}) \in \Lambda(\mathbf{c})$, where $\Lambda(\mathbf{c})$ is the set of optimal Lagrange multipliers associated with constraint (2.2).

4.3.1 Solution of the Multipath Rate Control Subproblem

Throughout the rest of this thesis, algorithms are presented using the continuous-time model. The presented continuous-time algorithms can also be viewed as the functional limit of their discrete-time counterparts, providing that the discrete-time time steps are appropriately rescaled and step sizes are close to zero.

Rate control problem (4.9) has been extensively studied in the literature [35]. Here, the duality-based solutions are considered where Lagrange variables are updated according to

$$\dot{\lambda}_l = \frac{\beta}{c_l} [R_l \mathbf{x}(\boldsymbol{\lambda}) - c_l]_{\lambda_l}^+ \quad \forall l \in L \quad (4.10)$$

where $\beta > 0$, R_l is the l th row of R , $\mathbf{x}(\boldsymbol{\lambda})$ are the path rates given $\boldsymbol{\lambda}$, and $[g(x)]_x^+$ is defined by

$$[g(x)]_x^+ = \begin{cases} g(x) & x > 0 \\ \max(g(x), 0) & x = 0 \end{cases}$$

Since the objective function in (4.9) is not strictly concave in \mathbf{x} , by Proposition 6.1.1 in [4], the dual function of (4.9) may not be differentiable at every point. Moreover, as shown in Section 6.1 in [4], the term within the brackets in (4.10) is a subgradient of the dual function and thus the term on right side of (4.10) is discontinuous. For the discrete-time version of (4.10), it is shown in [28] that while duality-based approaches always converge to a dual optimal solution, as a result of non-differentiability of the dual function, path rates \mathbf{x} do not converge and continuously oscillate.

For certain forms of R , such as the case where the number of disjoint paths is sufficiently large, optimal path rates $\mathbf{x}(\mathbf{c})$ are unique, as shown in the following lemma.

Lemma 4.1. *Let I_d be the number of disjoint paths in R , i.e. paths that do not share any links with any other paths. Furthermore, let S_d be the number of sources with only disjoint paths. If*

$$I_d = I - S + S_d \quad (4.11)$$

then $\mathbf{x}(\mathbf{c})$ is the unique primal optimal solution of (4.9).

Proof. Let $B_L = \{l \in L | R_l \mathbf{x}(\mathbf{c}) = c_l\}$ and $B_I = \{i \in I_s, s \in S | x_i^s(\mathbf{c}) = 0\}$. It can be easily verified that at optimality, for every $i \notin B_I$, there exists $l \in B_L$ such that $R_{l,i}^s = 1$ and $\lambda_l(\mathbf{c}) > 0$. The strict complementary slackness condition at l results from (4.4) and the assumption $f' > 0$. Moreover, if $i \in I_d$ then $i \notin B_I$. Since disjoint paths do not share any links, it then follows that there exist at least I_d linearly independent vectors R_l , where $l \in B_L$ and $\lambda_l(\mathbf{c}) > 0$. Further, it results from the definition of Lagrangian in (4.3) that

$$\begin{aligned} \mathbf{z}^T \nabla_{\mathbf{x}}^2 L(\mathbf{x}(\mathbf{c}), \boldsymbol{\lambda}(\mathbf{c}), \boldsymbol{\mu}(\mathbf{c})) \mathbf{z} &= \sum_{s \in S} w_s f_s''(x_s(\mathbf{c})) \left(\sum_{i \in I_s} z_i \right)^2 \\ &\leq 0 \quad \forall \mathbf{z} \neq \mathbf{0} \end{aligned} \quad (4.12)$$

Clearly, in order for (4.12) to be zero, $\sum_{i \in I_s} z_i = 0$, for all $s \in S$. This results in S linear equations of \mathbf{z} with linearly independent multiplier vectors, which, except for sources with only disjoint paths, are also linearly independent from R_l vectors, where $l \in B_L$ and $\lambda_l(\mathbf{c}) > 0$, that are associated with $i \in I_d$. Hence, given (4.11), there does not exist $\mathbf{z} \in \mathbb{R}^I$, $\mathbf{z} \neq \mathbf{0}$ such that (4.12) is zero and $R_l \mathbf{z} = 0$, for all $l \in B_L$, $\lambda_l(\mathbf{c}) > 0$. It then follows that the second-order sufficient conditions for a unique local maximising point of (4.9) (Lemma 3.2.1 in [14], Appendix B.1) hold at $\mathbf{x}(\mathbf{c})$. Furthermore, concavity of (4.9) implies that $\mathbf{x}(\mathbf{c})$ is the unique global maximiser of (4.9). \square

The right-hand side of algorithm (4.10) corresponds to the β multiple of marginal increase in average queueing delay at link l , given the traffic rate entering link l is equal to $R_l \mathbf{x}(\boldsymbol{\lambda})$. While this condition holds at links at the network traffic entry points, the traffic rate at the other links are bounded by the data rates of the links connected to their source node. However, at the equilibrium, the right-hand side of (4.10) is the β multiple of marginal increase in average queueing delay for all links $l \in L$. This implies that, by the results from stability of systems with vanishing perturbation [21], if link prices are updated according to β multiple of the link average queueing delays, path rates $\mathbf{x}(\boldsymbol{\lambda})$ and link prices $\boldsymbol{\lambda}$ converge to the primal and dual optimal solutions of (4.9), respectively, given β is sufficiently small. In this case, link average queueing delays at equilibrium are equal to $\beta^{-1} \boldsymbol{\lambda}(\mathbf{c})$.

4.3.2 Solution of the Scheduling Problem

First, the condition under which optimal link data rates are unique is given in the following lemma.

Lemma 4.2. *Let \mathbf{c}^* and $\tilde{\mathbf{c}}^*$ be optimal link data rates for (4.1). Let $\mathbf{x}(\mathbf{c}^*)$ and $\mathbf{x}(\tilde{\mathbf{c}}^*)$ be optimal path rates corresponding to \mathbf{c}^* and $\tilde{\mathbf{c}}^*$, respectively. Then $\mathbf{x}(\mathbf{c}^*) = \mathbf{x}(\tilde{\mathbf{c}}^*)$ implies $\mathbf{c}^* = \tilde{\mathbf{c}}^*$.*

Proof. Suppose that $\mathbf{c}^* \neq \tilde{\mathbf{c}}^*$. Since the objective function of (4.3) is an affine function of \mathbf{c} , and \mathbf{c}^* and $\tilde{\mathbf{c}}^*$ are maximisers of (4.3) at $(\boldsymbol{\lambda}^*, \boldsymbol{\mu}^*)$, they are maximal points of $\text{Co}(C)$, i.e. if $\mathbf{c} \in \text{Co}(C)$, $\mathbf{c} \succeq \mathbf{c}^*[\tilde{\mathbf{c}}^*]$ only if $\mathbf{c} = \mathbf{c}^*[\tilde{\mathbf{c}}^*]$. Thus, there exist $l, \hat{l} \in L$, $l \neq \hat{l}$ such that $c_l^* < \tilde{c}_l^*$ and $c_{\hat{l}}^* > \tilde{c}_{\hat{l}}^*$. Let $\mathbf{x}(\mathbf{c}^*) = \mathbf{x}(\tilde{\mathbf{c}}^*) \triangleq \mathbf{x}^*$. It then follows from (2.2) that $R_l \mathbf{x}^* \leq c_l^* < \tilde{c}_l^*$ and $R_{\hat{l}} \mathbf{x}^* \leq \tilde{c}_{\hat{l}}^* < c_{\hat{l}}^*$, which means that $R\mathbf{x}^*$ is upper bounded by a link data rates vector that is not a maximal point of $\text{Co}(C)$. This implies that \mathbf{x}^* is not optimal, which contradicts the initial assumption. \square

Inspired by the ideas from the gradient optimisation methods [4], the following solution for scheduling problem (4.8) is proposed

$$\dot{\mathbf{c}} = \gamma(\tilde{\mathbf{c}} - \mathbf{c}) \quad (4.13)$$

where $0 < \gamma \leq 1$, and

$$\tilde{\mathbf{c}} = \begin{cases} \mathbf{c} & \mathbf{c} = \arg \max_{\boldsymbol{\varsigma} \in \text{Co}(C)} \boldsymbol{\lambda}(\mathbf{c})^T \boldsymbol{\varsigma} \\ \arg \max_{\boldsymbol{\varsigma} \in C} \boldsymbol{\lambda}(\mathbf{c})^T \boldsymbol{\varsigma} & \text{otherwise} \end{cases} \quad (4.14)$$

where $\boldsymbol{\lambda}(\mathbf{c})$ is the optimal Lagrange variable of (4.9) given \mathbf{c} . Since C is a finite set, $\text{Co}(C)$ is a polyhedral set and hence by Proposition B.21 in [4], the optimisation problem in (4.14) attains a maximum at some extreme point of $\text{Co}(C)$. Therefore, the solution space in (4.14) is reduced to C . It can be seen that $\boldsymbol{\lambda}(\mathbf{c})^T (\tilde{\mathbf{c}} - \mathbf{c}) > 0$ when \mathbf{c} is not an equilibrium. Given $\boldsymbol{\lambda}(\mathbf{c})$ is a subgradient of U_w at \mathbf{c} , then (4.13)-(4.14) resembles a gradient optimisation method. The problem (4.14) is of the same form as the well-known scheduling problem (3.10) and thus can be solved using distributed solutions discussed in Section 3.3.

The right-hand side of (4.13) may not be continuous in general, for example, when $\lambda(\mathbf{c})$ is not unique, or when in (4.14) strict complimentary slackness condition does not hold at $\tilde{\mathbf{c}}$ (Theorem 3.2.2 in [14], Appendix B.1), and as a result the existence of solutions is not guaranteed. In the analysis that follows it is assumed that the following conditions, which are prerequisites for the results in [2], hold

H1 For any initial condition \mathbf{c}_0 , at least one solution of (4.13)-(4.14) exists.

H2 The right-hand side of (4.13) is Lebesgue measurable and locally bounded.

The following definitions apply Definitions 3, 4, and 6 in [2] (Appendix B.2) to algorithm (4.13)-(4.14):

Definition 1. A function $V : \mathbf{R}^L \rightarrow \mathbf{R}$ is said to be *nonpathological* if it is locally Lipschitz continuous and for every absolutely continuous function $\mathbf{c} : T \rightarrow \mathbf{R}^L$ and for almost every $t \in T$, the set $\partial_C V(\mathbf{c}(t))$ is a subset of an affine subspace orthogonal to $\dot{\mathbf{c}}(t)$, where $\partial_C V(\mathbf{c})$ denotes the Clarke gradient of real function V at point \mathbf{c} .

Definition 2. Let $V : \mathbf{R}^L \rightarrow \mathbf{R}$ be a nonpathological function and $\mathbf{g}(\mathbf{c})$ denote the right-hand side of (4.13). Let

$$A_V = \left\{ \mathbf{c} \in \mathbf{R}^L : \mathbf{e}_1^T \mathbf{g}(\mathbf{c}) = \mathbf{e}_2^T \mathbf{g}(\mathbf{c}) \quad \forall \mathbf{e}_1, \mathbf{e}_2 \in \partial_C V(\mathbf{c}) \right\} \quad (4.15)$$

if $\mathbf{c} \in A_V$, the *nonpathological derivative of the map V with respect to (4.13)-(4.14) at \mathbf{c}* is defined by

$$\dot{\bar{V}}_{\mathbf{g}}(\mathbf{c}) = \mathbf{e}^T \mathbf{g}(\mathbf{c}) \quad (4.16)$$

where \mathbf{e} is any vector in $\partial_C V(\mathbf{c})$.

Definition 3. A set M is said to be *weakly invariant* for (4.13)-(4.14) if for any $\mathbf{c}_0 \in M$ there exists a $\mathbf{c} \in S_{\mathbf{c}_0}$, where $S_{\mathbf{c}_0}$ denotes the set of maximal solutions of (4.13)-(4.14) with initial condition \mathbf{c}_0 , such that $\mathbf{c}(t) \in M$ for all $t \geq 0$.

By Theorem 2.2.6 in [14] (Appendix B.1), the mapping $\mathbf{q}(\mathbf{c})$ is continuous, and since $\mathbf{q}(\mathbf{c})$ is also unique, it is a continuous function. It is assumed that $\mathbf{q}(\mathbf{c})$ is also nonpathological [2].

Theorem 4.1. *Algorithms (4.13)-(4.14) converge to an optimal solution of (4.8).*

Proof. Let $\hat{\mathbf{c}}$ be an equilibrium point of (4.13)-(4.14). From (4.14) it follows that

$$\hat{\mathbf{c}} = \arg \max_{\mathbf{c} \in \text{Co}(C)} \boldsymbol{\lambda}(\hat{\mathbf{c}})^T \boldsymbol{\varsigma} \quad (4.17)$$

Since $\boldsymbol{\lambda}(\hat{\mathbf{c}})$ is a subgradient of U_w at $\hat{\mathbf{c}}$,

$$U_w(\mathbf{c}) \leq U_w(\hat{\mathbf{c}}) + \boldsymbol{\lambda}(\hat{\mathbf{c}})^T (\mathbf{c} - \hat{\mathbf{c}}) \quad \forall \mathbf{c} \in \text{Co}(C)$$

It follows from (4.17) that $\boldsymbol{\lambda}(\hat{\mathbf{c}})^T (\mathbf{c} - \hat{\mathbf{c}}) \leq 0$, so $U_w(\mathbf{c}) \leq U_w(\hat{\mathbf{c}})$ for all $\mathbf{c} \in \text{Co}(C)$. This means that $\hat{\mathbf{c}}$ is an optimal solution of (4.8), i.e. $\hat{\mathbf{c}} \in C^*$, where C^* denotes the set of optimal solutions of (4.8).

Consider the Lyapunov function

$$V(\mathbf{c}) = \frac{1}{2} \|\mathbf{q}(\mathbf{c}^*) - \mathbf{q}(\mathbf{c})\|_2^2$$

where $\mathbf{c}^* \in C^*$. Since $\mathbf{q}(\mathbf{c}^*) = \mathbf{q}^*$ is unique, $V(\mathbf{c}^*) = 0$ and $V(\mathbf{c}) > 0$, for all $\mathbf{c} \notin C^*$. Moreover, since $\mathbf{q}(\mathbf{c})$ is nonpathological, $V(\mathbf{c})$ is also nonpathological. Let $\dot{\bar{V}}$ be the nonpathological derivative of the map V with respect to (4.13)-(4.14) at $\mathbf{c} \in A_V$, where A_V and $\dot{\bar{V}}$ are defined in (4.16) and (4.15), respectively. Let $\boldsymbol{\psi}_s \in \partial_C q_s(\mathbf{c})$, $s \in S$, where $\partial_C q_s(\mathbf{c})$ is the Clarke gradient of q_s at \mathbf{c} . Also let $\Psi = [\boldsymbol{\psi}_1 \cdots \boldsymbol{\psi}_S]^T$. Then

$$\begin{aligned} \dot{\bar{V}}(\mathbf{c}) &= -(\mathbf{q}(\mathbf{c}^*) - \mathbf{q}(\mathbf{c}))^T \dot{\bar{\mathbf{q}}}(\mathbf{c}) \\ &= -(\mathbf{q}(\mathbf{c}^*) - \mathbf{q}(\mathbf{c}))^T \Psi \dot{\mathbf{c}} \\ &= -(\mathbf{q}(\mathbf{c}^*) - \mathbf{q}(\mathbf{c}))^T \Psi \gamma(\tilde{\mathbf{c}} - \mathbf{c}) \\ &= -\gamma(\mathbf{q}(\mathbf{c}^*) - \mathbf{q}(\mathbf{c}))^T \Psi (\mathbf{c}^* - \mathbf{c} + \tilde{\mathbf{c}} - \mathbf{c}^*) \\ &= -\gamma(\mathbf{q}(\mathbf{c}^*) - \mathbf{q}(\mathbf{c}))^T \Psi (\mathbf{c}^* - \mathbf{c}) \\ &\quad - \gamma(\mathbf{q}(\mathbf{c}^*) - \mathbf{q}(\mathbf{c}))^T \Psi (\tilde{\mathbf{c}} - \mathbf{c}^*) \end{aligned}$$

Using the characterisation of Clarke gradient in equation A.11 in [1] (Appendix B.2), it follows from Taylor's theorem that $\mathbf{q}(\mathbf{c}^*) - \mathbf{q}(\mathbf{c}) \approx \Psi(\mathbf{c}^* - \mathbf{c})$, as \mathbf{c} approaches \mathbf{c}^* . Furthermore, since U_w is concave, by Proposition B.24 in [4], there exists $\hat{\boldsymbol{\lambda}} \in \Lambda(\mathbf{c}^*)$ such that

$$\hat{\boldsymbol{\lambda}}^T (\mathbf{c} - \mathbf{c}^*) \leq 0 \quad \forall \mathbf{c} \in \text{Co}(C) \quad (4.18)$$

thus

$$\dot{\bar{V}}(\mathbf{c}) \approx -\gamma \|\mathbf{q}(\mathbf{c}^*) - \mathbf{q}(\mathbf{c})\|_2^2 - \gamma(\hat{\boldsymbol{\lambda}} - \boldsymbol{\lambda}(\mathbf{c}))^T R \Psi(\tilde{\mathbf{c}} - \mathbf{c}^*)$$

where $\hat{\boldsymbol{\lambda}}$ satisfies (4.18).

It can be shown that typically $R \Psi \approx -k(\mathbf{c})I_L$, where $k(\mathbf{c}) > 0$ and I_L is the identity matrix. To see this, consider the case where each source has only a single path, i.e. $I_s = 1$. In this case $\mathbf{x}(\boldsymbol{\lambda})$ is differentiable with respect to $\boldsymbol{\lambda}$ and $\frac{\partial \mathbf{x}(\boldsymbol{\lambda})}{\partial \boldsymbol{\lambda}} = \text{diag} \left\{ \frac{1}{w_s f_s''(x_s(\boldsymbol{\lambda}))} \right\} R^T$ [30]. Evaluating the sensitivity equation (2.9) in [21] for dual algorithm (4.10) at its equilibrium point $\boldsymbol{\lambda}(\mathbf{c})$ yields

$$\begin{aligned} 0 &= \text{diag} \left\{ \frac{\beta_l}{c_l} \right\} R \frac{\partial \mathbf{x}(\boldsymbol{\lambda}(\mathbf{c}))}{\partial \boldsymbol{\lambda}} \frac{\partial \boldsymbol{\lambda}(\mathbf{c})}{\partial \mathbf{c}} - \text{diag} \left\{ \frac{\beta_l}{c_l^2} R_l \mathbf{x}(\boldsymbol{\lambda}(\mathbf{c})) \right\} \\ &= \text{diag} \left\{ \frac{\beta_l}{c_l} \right\} R \text{diag} \left\{ \frac{1}{w_s f_s''(x_s(\boldsymbol{\lambda}(\mathbf{c})))} \right\} R^T \frac{\partial \boldsymbol{\lambda}(\mathbf{c})}{\partial \mathbf{c}} - \text{diag} \left\{ \frac{\beta_l}{c_l^2} R_l \mathbf{x}(\boldsymbol{\lambda}(\mathbf{c})) \right\} \\ &\approx R \text{diag} \left\{ \frac{1}{w_s f_s''(x_s(\boldsymbol{\lambda}(\mathbf{c})))} \right\} \frac{\partial \mathbf{q}(\mathbf{c})}{\partial \mathbf{c}} - I_L \end{aligned}$$

The last approximation is based on the fact that at optimality total flow on each link is near or equal its capacity. Thus, after factoring out $\text{diag} \left\{ \frac{\beta_l}{c_l} \right\}$, the second term on the right-hand side of the equality can be approximated as an identity matrix. Furthermore, it is assumed that the system operates at points where w_s and as a result x_s^* are close for all $s \in S$. Consequently, the values of $w_s f_s''(x_s(\mathbf{c}))$, $s \in S$ are close. Hence, $R \frac{\partial \mathbf{q}(\mathbf{c})}{\partial \mathbf{c}} \approx -k(\mathbf{c})I_L$, where $k(\mathbf{c}) \approx |w_s f_s''(x_s(\mathbf{c}))|$, $s \in S$.

From (4.14) it follows that

$$\boldsymbol{\lambda}(\mathbf{c})^T (\tilde{\mathbf{c}} - \mathbf{c}^*) \geq 0 \quad \forall \mathbf{c} \in \text{Co}(C)$$

Also, (4.18) implies

$$\hat{\boldsymbol{\lambda}}^T (\tilde{\mathbf{c}} - \mathbf{c}^*) \leq 0$$

Adding both inequalities yields

$$(\hat{\boldsymbol{\lambda}} - \boldsymbol{\lambda}(\mathbf{c}))^T (\tilde{\mathbf{c}} - \mathbf{c}^*) \leq 0 \quad \forall \mathbf{c} \in \text{Co}(C) \quad (4.19)$$

thus

$$\begin{aligned} \dot{\bar{V}}(\mathbf{c}) &= -\gamma \|\mathbf{q}(\mathbf{c}^*) - \mathbf{q}(\mathbf{c})\|_2^2 + \gamma k(\mathbf{c}) I_L (\hat{\boldsymbol{\lambda}} - \boldsymbol{\lambda}(\mathbf{c}))^T (\tilde{\mathbf{c}} - \mathbf{c}^*) \\ &\leq -\gamma \|\mathbf{q}(\mathbf{c}^*) - \mathbf{q}(\mathbf{c})\|_2^2 \end{aligned} \quad (4.20)$$

Furthermore, the largest weakly invariant set of points $\mathbf{c} \in A_v$ for which $\dot{\bar{V}}(\mathbf{c}) = 0$ is the set of equilibrium points C^* . Hence, by Proposition 3 in [2] (Appendix B.2), every solution of algorithms (4.13)-(4.14) approaches the set of equilibrium points C^* as $t \rightarrow \infty$. \square

4.4 Conclusions

In this chapter the optimisation problem (4.1) is presented as an alternative formulation to the original optimisation problem (2.1)-(2.5). In the alternative formulation the delay constraints are omitted and utility functions are multiplied by weight coefficients. Generally, primal optimal solutions $(\mathbf{x}^*, \mathbf{c}^*)$ and dual optimal solutions, or optimal link prices, $\boldsymbol{\lambda}^*$ are not unique. However, optimal aggregate source rates $\{x_s^*\}$ are unique. Furthermore, paths with positive flows at optimality have minimum and thus equal price. It then follows that optimal path prices are \mathbf{q}^* are also unique.

The alternative optimisation problem (4.1) is then presented as an equivalent scheduling problem, which consists of the master scheduling problem (4.8) and the well-known multi-path rate control with fixed link data rates subproblem (4.9). The rate control subproblem is solved using the duality-based algorithm, where dual variables, or link prices, are updated in proportion to the link queueing delays. Using this update mechanism, link queueing delays at equilibrium are proportional to the optimal link prices $\boldsymbol{\lambda}(\mathbf{c})$. However, it has been shown that while the duality-based approach always converge to a dual optimal solution, for sources with multiple alternative paths, path rates do not converge and continuously oscillate.

Next, conditions on the number of disjoint paths are derived that guarantee unique optimal path rates $\mathbf{x}^*(\mathbf{c})$. This result is based on satisfying the second-order sufficient conditions for a unique local maximising point of the multi-path rate control subproblem (4.9). This suggests an approach for future work where the second-order sufficient conditions are used to design a multi-path rate control algorithm that converges the unique optimal path rates, given the conditions on

the number of disjoint paths are guaranteed by the topology control algorithm.

A distributed algorithm for the scheduling problem (4.8) is then proposed and is shown to converge to an optimal solution \mathbf{c}^* . The proposed scheduling algorithm is based on solving the well-known scheduling problem (3.10) and thus can incorporate the distributed solutions discussed in Section 3.3.

The proposed solution for providing bounded delay comprises an algorithm that is incorporated in the proposed scheduling algorithm and is based on its key feature which is the proportionality of link queueing delays and link prices $\lambda(\mathbf{c})$, as well as the correlation between optimal path prices and utility weights coefficients in the alternative optimisation problem (4.1), as will be described in the next chapter.

Chapter 5

Proposed Solution for Providing Bounded Delay

5.1 Introduction

In this chapter the proposed solution to the optimisation problem (2.1)-(2.5) is developed, with focus on the performance objectives described in Section 1.3.4, namely, ensuring bounded end-to-end queueing delays, enabling distributed implementation with low communication overhead, leading to maximal link capacity utilisation, and having controllable transient behaviour. The proposed solution exploits the properties of the alternative optimisation problem (4.1) and its proposed solution described in Chapter 4. In Section 5.2, bounds on the sensitivity of optimal path prices $\mathbf{q}^*(\mathbf{w})$ and aggregate source rates to the variations of utility weight coefficients \mathbf{w} in the alternative optimisation problem (4.1) are derived. Based on the sensitivity results, in Section 5.3 an algorithm for providing bounded end-to-end queueing delays as well as other performance objectives is proposed, which is integrated in the scheduling algorithms (4.13)-(4.14).

5.2 Effect of Sources' Weights on Delay

As explained in Section 4.3.1, if link prices in the duality-based algorithm (4.10) are instead updated proportionally to link average queueing delays, optimal link prices are then proportional to the link average queueing delays at equilibrium. Furthermore, by (4.7) optimal path prices for each source $s \in S$ are equal to its marginal utility at its optimal aggregate data transmission rate, multiplied by its weight. This suggests an alternative approach to the original formulation (2.1)-(2.5), in which the delay bounds in (2.4) are instead guaranteed by adjusting the weight of sources. Optimal allocation of sources' data transmission rates in this approach may differ slightly from the case with the initial source weights but no delay constraints, and consequently may lead to slight reduction in the overall perceived signal quality. However, since this approach utilises the network available capacity, it outperforms the previously proposed solutions described in Section 3 in terms of the overall perceived signal quality.

The following lemma shows that for the alternative problem (4.1) the optimal path prices for each source $s \in S$, that is $q_s^* = R_i^{sT} \lambda^*$, $i \in I_s$, grow as its weight w_s increases.

Lemma 5.1. *If the utility functions f_s are twice continuously differentiable, strictly concave, increasing, and $f_s'' < 0$ for all $s \in S$, then upper and lower bounds on the sensitivity of $q_s^*(w)$ and $x_s^*(w)$ to the variation of parameters w_s are given by*

$$0 < \frac{\partial q_s^*}{\partial w_s} \leq f'_s(x_s^*) \quad (5.1)$$

$$0 \leq \frac{\partial x_s^*}{\partial w_s} < \frac{f'_s(x_s^*)}{w_s f''_s(x_s^*)} \quad (5.2)$$

for all $s \in S$.

Proof. It results from (4.7) that

$$\frac{\partial q_s^*}{\partial w_r} = \begin{cases} w_s f''_s(x_s^*) \frac{\partial x_s^*}{\partial w_s} + f'_s(x_s^*) & r = s \\ w_s f''_s(x_s^*) \frac{\partial x_s^*}{\partial w_r} & r \neq s \end{cases} \quad \forall s \in S \quad (5.3)$$

Let $\tilde{\mathbf{w}}$ be a perturbation of parameter \mathbf{w} defined by

$$\tilde{w}_s = \begin{cases} w_s + dw_r & s = r \\ w_s & \text{otherwise} \end{cases} \quad \forall s \in S$$

where $r \in S$ and $dw_r > 0$. If $x_r^*(\tilde{\mathbf{w}}) = x_r^*(\mathbf{w})$, taking the limit $dw_r \rightarrow 0$ yields $\frac{\partial x_r^*}{\partial w_r} = 0$, and from (5.3), $\frac{\partial q_r^*}{\partial w_r} = f'_r(x_r^*)$. If $x_r^*(\tilde{\mathbf{w}}) \neq x_r^*(\mathbf{w})$, given the strict concavity of f , $\{x_s^*(\mathbf{w})\}$ and $\{x_s^*(\tilde{\mathbf{w}})\}$ are the unique maximisers for problem (4.1) with parameters \mathbf{w} and $\tilde{\mathbf{w}}$, respectively. So

$$\sum_{s \in S} w_s f_s(x_s^*(\mathbf{w})) > \sum_{s \in S} w_s f_s(x_s^*(\tilde{\mathbf{w}}))$$

and

$$\sum_{s \in S} \tilde{w}_s f_s(x_s^*(\tilde{\mathbf{w}})) > \sum_{s \in S} \tilde{w}_s f_s(x_s^*(\mathbf{w}))$$

Adding both inequalities results in

$$\sum_{s \in S} (\tilde{w}_s - w_s) (f_s(x_s^*(\tilde{\mathbf{w}})) - f_s(x_s^*(\mathbf{w}))) > 0$$

Except for $s = r$, all the elements in the above summation are zero. Since $\tilde{w}_r - w_r = dw_r > 0$, $f_r(x_r^*(\tilde{\mathbf{w}})) > f_r(x_r^*(\mathbf{w}))$, which implies $x_r^*(\tilde{\mathbf{w}}) > x_r^*(\mathbf{w})$, since f is an increasing function. Thus

$$\frac{x_r^*(\tilde{\mathbf{w}}) - x_r^*(\mathbf{w})}{dw_r} > 0$$

Taking the limit $dw_r \rightarrow 0$ yields the lower bound of (5.2).

It results from the optimality condition in Proposition 2.2.2 in [4] for optimisation problem (4.1) at \mathbf{w} that

$$\sum_{s \in S} w_s f'_s(x_s^*(\mathbf{w})) (x_s^*(\tilde{\mathbf{w}}) - x_s^*(\mathbf{w})) \leq 0$$

Similarly, it results from the optimality condition for (4.1) at the perturbed $\tilde{\mathbf{w}}$ that

$$\sum_{s \in S} \tilde{w}_s f'_s(x_s^*(\tilde{\mathbf{w}})) (x_s^*(\mathbf{w}) - x_s^*(\tilde{\mathbf{w}})) \leq 0$$

Using definition (4.7), adding both inequalities and taking the limit $dw_r \rightarrow 0$ yields

$$\sum_{s \in S} \frac{\partial x_s^*}{\partial w_r} \frac{\partial q_s^*}{\partial w_r} \geq 0 \quad \forall r \in S \quad (5.4)$$

Let $S_d = \{s \in S, s \neq r | x_s^*(\tilde{\mathbf{w}}) < x_s^*(\mathbf{w})\}$. Since $x_r^*(\tilde{\mathbf{w}}) > x_r^*(\mathbf{w})$, S_d is non-empty, otherwise $\{x_s^*(\mathbf{w})\}$ would not be optimal. Since it is assumed that $f'' < 0$, $f'_s(x_s^*(\tilde{\mathbf{w}})) > f'_s(x_s^*(\mathbf{w}))$, and hence from (4.7), $q_s^*(\tilde{\mathbf{w}}) > q_s^*(\mathbf{w})$ for all $s \in S_d$. Taking the limit $dw_r \rightarrow 0$ results in $\frac{\partial x_s^*}{\partial w_r} < 0$ and $\frac{\partial q_s^*}{\partial w_r} > 0$, so

$$\sum_{s \in S_d} \frac{\partial x_s^*}{\partial w_r} \frac{\partial q_s^*}{\partial w_r} < 0 \quad \forall r \in S$$

Let $S_i = \{s \in S, s \neq r | x_s^*(\tilde{\mathbf{w}}) \geq x_s^*(\mathbf{w})\}$. Using a similar argument, $q_s^*(\tilde{\mathbf{w}}) \leq q_s^*(\mathbf{w})$, for all $s \in S_i$, so

$$\sum_{s \in S_i} \frac{\partial x_s^*}{\partial w_r} \frac{\partial q_s^*}{\partial w_r} \leq 0 \quad \forall r \in S$$

Hence, since it was assumed that $\frac{\partial x_s^*}{\partial w_s} > 0$, it follows from (5.4) that $\frac{\partial q_s^*}{\partial w_s} > 0$, for all $s \in S$.

In (5.3), since $f''_s(x_s^*) < 0$, the first and second terms on the right side of the equation are negative and positive, respectively. Since the term on the left side of the equation is positive,

$$0 \leq -w_s f''_s(x_s^*) \frac{\partial x_s^*}{\partial w_s} < f'_s(x_s^*)$$

from which the upper bounds in (5.1) and (5.2) can be verified. \square

As explained in Section 4.3.1, if link prices in the duality-based rate control algorithm are updated according to β multiple of the link average queueing delays, link average queueing delays at equilibrium are equal to $\beta^{-1} \boldsymbol{\lambda}(\mathbf{c})$, and as a result path average queueing delays at equilibrium equal $\beta^{-1} \mathbf{q}(\mathbf{c})$. Consequently, Lemma (5.1) suggests a strategy for providing bounded end-to-end delay based on the adjustment of sources' weights.

5.3 Delay Regulation via Dynamic Adjustment of Sources' Weights

The main challenge in guaranteeing bounded delay through adjustment of sources' weights is that sources' weights that guarantee the required bounded delay generally vary for different network configurations. Clearly, a concave utility function

always implies that a connection has elastic bandwidth requirements and as a result best-effort strategies may allocate different rates to the same connection (with the same weight) under various network configurations in order to maximise the aggregate utility of all connections. This means that for every network configuration, sources' weights have to be recomputed to ensure the required bounded delay. Hence, the dynamic and decentralised nature of multihop wireless networks calls for a robust, responsive and distributed algorithm that can adjust sources' weights so as to ensure bounded end-to-end delay under modest parameter perturbations.

To order to meet these requirements, the following integral controller is proposed where, based on the results from Lemma (5.1), each source uses current end-to-end delay on its paths to adjust its weight w_s , and hence to regulate end-to-end delay at optimality.

$$\dot{w}_s = \alpha \left[d_s - \frac{q_s(\mathbf{c}, \mathbf{w})}{\beta} \right]_{w_s}^+ \quad \forall s \in S \quad (5.5)$$

Algorithm (5.5) is performed by each source independently and ensures bounded end-to-end delay under parameter perturbations that do not destabilise the system.

As in Section 4.3.2, in the following analysis it is assumed that the following conditions, which are prerequisites for the results in [2], hold

- H1** For any initial condition $(\mathbf{c}_0, \mathbf{w}_0)$, at least one solution of (4.13)-(4.14) and (5.5) exists.
- H2** The right-hand sides of (4.13) and (5.5) is Lebesgue measurable and locally bounded.

The following definitions apply Definitions 3, 4 and 6 in [2] (Appendix B.2) to algorithm (4.13)-(4.14) and (5.5):

Definition 1. A function $V : \mathbf{R}^{L+S} \rightarrow \mathbf{R}$ is said to be *nonpathological* if it is locally Lipschitz continuous and for every absolutely continuous function $(\mathbf{c}, \mathbf{w}) : T \rightarrow \mathbf{R}^{L+S}$ and for almost every $t \in T$, the set $\partial_C V(\mathbf{c}(t), \mathbf{w}(t))$ is a subset of an affine subspace orthogonal to $(\dot{\mathbf{c}}(t), \dot{\mathbf{w}}(t))$, where $\partial_C V(\mathbf{c}, \mathbf{w})$ denotes the Clarke gradient of real function V at point (\mathbf{c}, \mathbf{w}) .

Definition 2. Let $V : \mathbf{R}^{L+S} \rightarrow \mathbf{R}$ be a nonpathological function and $\mathbf{g}(\mathbf{c}, \mathbf{w})$ denote the right-hand side of (4.13) and (5.5). Let

$$A_V = \left\{ (\mathbf{c}, \mathbf{w}) \in \mathbf{R}^{L+S} : \mathbf{e}_1^T \mathbf{g}(\mathbf{c}, \mathbf{w}) = \mathbf{e}_2^T \mathbf{g}(\mathbf{c}, \mathbf{w}) \quad \forall \mathbf{e}_1, \mathbf{e}_2 \in \partial_C V(\mathbf{c}, \mathbf{w}) \right\} \quad (5.6)$$

if $(\mathbf{c}, \mathbf{w}) \in A_V$, the nonpathological derivative of the map V with respect to (4.13)-(4.14) and (5.5) at (\mathbf{c}, \mathbf{w}) is defined by

$$\dot{\bar{V}}_g(\mathbf{c}, \mathbf{w}) = \mathbf{e}^T \mathbf{g}(\mathbf{c}, \mathbf{w}) \quad (5.7)$$

where \mathbf{e} is any vector in $\partial_C V(\mathbf{c}, \mathbf{w})$.

Definition 3. A set M is said to be *weakly invariant* for (4.13)-(4.14) and (5.5) if for any $(\mathbf{c}_0, \mathbf{w}_0) \in M$ there exists a $(\mathbf{c}, \mathbf{w}) \in S_{(\mathbf{c}_0, \mathbf{w}_0)}$, where $S_{(\mathbf{c}_0, \mathbf{w}_0)}$ denotes the set of maximal solutions of (4.13)-(4.14) and (5.5) with initial condition $(\mathbf{c}_0, \mathbf{w}_0)$, such that $(\mathbf{c}(t), \mathbf{w}(t)) \in M$ for all $t \geq 0$.

The following theorem examines the conditions under which algorithms (4.13)-(4.14) combined with (5.5) achieve asymptotic regulation of end-to-end delay.

Theorem 5.1. *Algorithms (4.13)-(4.14) combined with (5.5) converge to an optimal solution of (4.8) with parameter \mathbf{w}^* , where \mathbf{w}^* is the weight of sources that guarantees bounded delay specified in (2.4), if*

- *subproblem (4.9) is solved using duality-based algorithm (4.10), where link prices are instead updated as β multiple of link average queueing delays,*
- *the following conditions hold*

$$\left| 1 - \frac{1}{\epsilon_s} \right| \leq \frac{1}{\sqrt{S-1}} \quad \forall s \in S \quad (5.8)$$

where $0 < \epsilon_s \leq 1$, $s \in S$ satisfy

$$\frac{\partial q_s^*}{\partial w_s} = \epsilon_s f'_s(x_s^*) \quad \forall s \in S \quad (5.9)$$

- *and parameter α in (5.5) satisfies*

$$\alpha < \frac{\beta \gamma}{2 \max_{s \in S} f'_s(x_s^*)} \quad (5.10)$$

Proof. As explained in Section 4.3.1, if link prices in the duality-based rate control algorithm are updated according to β multiple of the link average queueing delays, link average queueing delays at equilibrium are equal to $\beta^{-1}\boldsymbol{\lambda}(\mathbf{c})$, and as a result path average queueing delays at equilibrium equal $\beta^{-1}\mathbf{q}(\mathbf{c})$, which are bounded by \mathbf{d} at the equilibrium $(\mathbf{c}^*, \mathbf{w}^*)$.

By Theorem 2.2.6 in [14] (Appendix B.1), the mapping $\mathbf{q}^*(\mathbf{w})$ is continuous, and since $\mathbf{q}^*(\mathbf{w})$ is also unique, it is a continuous function. It is assumed that $\mathbf{q}^*(\mathbf{w})$ is also nonpathological. Consider the Lyapunov function

$$V(\mathbf{c}, \mathbf{w}) = \frac{1}{2} \|\mathbf{q}^*(\mathbf{w}^*) - \mathbf{q}^*(\mathbf{w})\|_2^2 + \frac{1}{2} \|\mathbf{q}^*(\mathbf{w}) - \mathbf{q}_w(\mathbf{c})\|_2^2$$

Where $\mathbf{q}_{w^*}(\mathbf{c}^*) = \mathbf{q}^*(\mathbf{w}^*) = \beta\mathbf{d}$. Therefore, $V(\mathbf{c}^*, \mathbf{w}^*) = 0$ and $V(\mathbf{c}, \mathbf{w}) > 0$, for all $(\mathbf{c}, \mathbf{w}) \neq (\mathbf{c}^*, \mathbf{w}^*)$. Moreover, since $\mathbf{q}_w(\mathbf{c})$ and $\mathbf{q}^*(\mathbf{w})$ are nonpathological, $V(\mathbf{c}, \mathbf{w})$ is also nonpathological. Let $\dot{\bar{V}}$ be the nonpathological derivative of the map V with respect to (4.13)-(4.14) and (5.5) at $(\mathbf{c}, \mathbf{w}) \in A_V$, where A_V and $\dot{\bar{V}}$ are defined in (5.6) and (5.7), respectively. Let $\Psi_w = [\psi_{1,w} \cdots \psi_{S,w}]^T$, where $\psi_{s,w} \in \partial_C q_{s,w}(\mathbf{c})$, $s \in S$, and $\partial_C q_{s,w}(\mathbf{c})$ is the Clarke gradient of $q_{s,w}$ at \mathbf{c} . Also let $\Phi = [\phi_1 \cdots \phi_S]^T$, where $\phi_s \in \partial_C q_s^*(\mathbf{w})$, $s \in S$, and $\partial_C q_s^*(\mathbf{w})$ is the Clarke gradient of q_s^* at \mathbf{w} . Then

$$\begin{aligned} \dot{\bar{V}}(\mathbf{c}, \mathbf{w}) &= -(\mathbf{q}^*(\mathbf{w}^*) - \mathbf{q}^*(\mathbf{w}))^T \dot{\bar{\mathbf{q}}}^*(\mathbf{w}) + (\mathbf{q}^*(\mathbf{w}) - \mathbf{q}_w(\mathbf{c}))^T (\dot{\bar{\mathbf{q}}}^*(\mathbf{w}) - \dot{\bar{\mathbf{q}}}_w(\mathbf{c})) \\ &= -(\mathbf{q}^*(\mathbf{w}^*) - \mathbf{q}^*(\mathbf{w}))^T \Phi \dot{w} + (\mathbf{q}^*(\mathbf{w}) - \mathbf{q}_w(\mathbf{c}))^T \Phi \dot{w} \\ &\quad - (\mathbf{q}^*(\mathbf{w}) - \mathbf{q}_w(\mathbf{c}))^T \dot{\bar{\mathbf{q}}}_w(\mathbf{c}) \\ &= -\frac{\alpha}{\beta} (\mathbf{q}^*(\mathbf{w}^*) - \mathbf{q}^*(\mathbf{w}))^T \Phi (\beta\mathbf{d} - \mathbf{q}_w(\mathbf{c})) \\ &\quad + \frac{\alpha}{\beta} (\mathbf{q}^*(\mathbf{w}) - \mathbf{q}_w(\mathbf{c}))^T \Phi (\beta\mathbf{d} - \mathbf{q}^*(\mathbf{w}) + \mathbf{q}^*(\mathbf{w}) - \mathbf{q}_w(\mathbf{c})) \\ &\quad - (\mathbf{q}^*(\mathbf{w}) - \mathbf{q}_w(\mathbf{c}))^T \dot{\bar{\mathbf{q}}}_w(\mathbf{c}) \\ &= -\frac{\alpha}{\beta} (\mathbf{q}^*(\mathbf{w}^*) - \mathbf{q}^*(\mathbf{w}))^T \Phi (\beta\mathbf{d} - \mathbf{q}_w(\mathbf{c})) \\ &\quad + \frac{\alpha}{\beta} (\mathbf{q}^*(\mathbf{w}) - \mathbf{q}_w(\mathbf{c}))^T \Phi (\beta\mathbf{d} - \mathbf{q}^*(\mathbf{w})) \\ &\quad + \frac{\alpha}{\beta} (\mathbf{q}^*(\mathbf{w}) - \mathbf{q}_w(\mathbf{c}))^T \Phi (\mathbf{q}^*(\mathbf{w}) - \mathbf{q}_w(\mathbf{c})) - (\mathbf{q}^*(\mathbf{w}) - \mathbf{q}_w(\mathbf{c}))^T \dot{\bar{\mathbf{q}}}_w(\mathbf{c}) \end{aligned}$$

Since $f'' < 0$, (5.3) implies that $\frac{\partial x_s^*}{\partial w_r} \frac{\partial q_s^*}{\partial w_r} \leq 0$, for all $r \neq s$. It then follows from

(5.4) and (5.3) that

$$\begin{aligned} \frac{\partial x_s^*}{\partial w_s} \frac{\partial q_s^*}{\partial w_s} &\geq \left| \sum_{\substack{r \in S \\ r \neq s}} \frac{\partial x_s^*}{\partial w_r} \frac{\partial q_s^*}{\partial w_r} \right| \quad \forall s \in S \\ \left(\frac{\partial q_s^*}{\partial w_s} - f'_s(x_s^*) \right) \frac{\partial q_s^*}{\partial w_s} \frac{1}{w_s f''_s(x_s^*)} &\geq \left| \sum_{\substack{r \in S \\ r \neq s}} \left(\frac{\partial q_s^*}{\partial w_r} \right)^2 \frac{1}{w_r f''_r(x_r^*)} \right| \quad \forall s \in S \end{aligned}$$

It is assumed that the system operates at points where w_s and as a result x_s^* are close for all $s \in S$. Therefore, the values of $w_s f''_s(x_s(\mathbf{c}))$, $s \in S$ are close. Furthermore, it follows from (5.1) that there exists $0 < \epsilon_s \leq 1$, $s \in S$ that satisfies (5.9) and so $\left| \frac{\partial q_s^*}{\partial w_s} - f'_s(x_s^*) \right| = \left| 1 - \frac{1}{\epsilon_s} \right| \frac{\partial q_s^*}{\partial w_s}$. Thus

$$\begin{aligned} \left| 1 - \frac{1}{\epsilon_s} \right| \frac{\partial q_s^*}{\partial w_s} &> \left(\sum_{\substack{r \in S \\ r \neq s}} \left(\frac{\partial q_s^*}{\partial w_r} \right)^2 \right)^{\frac{1}{2}} \\ &\geq \frac{1}{\sqrt{S-1}} \sum_{\substack{r \in S \\ r \neq s}} \left| \frac{\partial q_s^*}{\partial w_r} \right| \end{aligned}$$

It then follows from condition (5.8) that

$$\frac{\partial q_s^*}{\partial w_s} \geq \sum_{\substack{r \in S \\ r \neq s}} \left| \frac{\partial q_s^*}{\partial w_r} \right| \quad (5.11)$$

Inequality (5.11) implies that Φ is approximately strictly diagonally dominant (Definition 6.1.9 in [18], Appendix B.3). Moreover, off-diagonal elements of Φ are very small relative to the diagonal elements, and as a result Φ can be assumed to have almost the same properties as a symmetric matrix. Since by (5.1) the diagonal elements of Φ are positive, it then follows from Theorem 6.1.10 in [18] (Appendix B.3) that all eigenvalues of Φ are real and positive and hence Φ is

positive definite. Thus

$$\begin{aligned}
 \dot{\bar{V}}(\mathbf{c}, \mathbf{w}) &= -\frac{\alpha}{\beta}(\mathbf{q}^*(\mathbf{w}^*) - \mathbf{q}^*(\mathbf{w}))^T \Phi (\beta \mathbf{d} - \mathbf{q}_w(\mathbf{c})) \\
 &\quad + \frac{\alpha}{\beta}(\beta \mathbf{d} - \mathbf{q}^*(\mathbf{w}))^T \Phi (\mathbf{q}^*(\mathbf{w}) - \mathbf{q}_w(\mathbf{c})) \\
 &\quad + \frac{\alpha}{\beta}(\mathbf{q}^*(\mathbf{w}) - \mathbf{q}_w(\mathbf{c}))^T \Phi (\mathbf{q}^*(\mathbf{w}) - \mathbf{q}_w(\mathbf{c})) - (\mathbf{q}^*(\mathbf{w}) - \mathbf{q}_w(\mathbf{c}))^T \dot{\bar{\mathbf{q}}}_w(\mathbf{c}) \\
 &= -\frac{\alpha}{\beta}(\mathbf{q}^*(\mathbf{w}^*) - \mathbf{q}^*(\mathbf{w}))^T \Phi (\mathbf{q}^*(\mathbf{w}^*) - \mathbf{q}^*(\mathbf{w})) \\
 &\quad + \frac{\alpha}{\beta}(\mathbf{q}^*(\mathbf{w}) - \mathbf{q}_w(\mathbf{c}))^T \Phi (\mathbf{q}^*(\mathbf{w}) - \mathbf{q}_w(\mathbf{c})) - (\mathbf{q}^*(\mathbf{w}) - \mathbf{q}_w(\mathbf{c}))^T \dot{\bar{\mathbf{q}}}_w(\mathbf{c}) \\
 &< \frac{\alpha}{\beta}(\mathbf{q}^*(\mathbf{w}) - \mathbf{q}_w(\mathbf{c}))^T \Phi (\mathbf{q}^*(\mathbf{w}) - \mathbf{q}_w(\mathbf{c})) - (\mathbf{q}^*(\mathbf{w}) - \mathbf{q}_w(\mathbf{c}))^T \dot{\bar{\mathbf{q}}}_w(\mathbf{c})
 \end{aligned}$$

In the first equality, the second term on the right-hand side results from the assumption that Φ is symmetric. The above inequality results from the positive definiteness of Φ . It follows from Geršgorin Theorem (Theorem 6.1.1 in [18], Appendix B.3) and inequalities (5.11) and (5.1) that eigenvalues of Φ are upper-bounded by $2f'_s(x_s^*)$, $s \in S$. Applying Rayleigh-Ritz Theorem (Theorem 4.2.2 in [18], Appendix B.3) to the first term, and using the upperbound (4.20) for the second term on the right-hand side of the above inequality then yields

$$\dot{\bar{V}}(\mathbf{c}, \mathbf{w}) < \left(\frac{2\alpha}{\beta} \max_{s \in S} f'_s(x_s^*) - \gamma \right) \|\mathbf{q}^*(\mathbf{w}) - \mathbf{q}_w(\mathbf{c})\|_2^2 \quad (5.12)$$

Consequently, if parameter α satisfies (5.10) then $\dot{\bar{V}}(\mathbf{c}, \mathbf{w}) \leq 0$ for all (\mathbf{c}, \mathbf{w}) . Furthermore, the largest weakly invariant set of points $(\mathbf{c}, \mathbf{w}) \in A_V$ for which $\dot{\bar{V}}(\mathbf{c}, \mathbf{w}) = 0$ is the set of equilibrium points $(\mathbf{c}^*, \mathbf{w}^*)$. Hence, by Proposition 3 in [2] (Appendix B.2), every solution of algorithms (4.13)-(4.14) combined with (5.5) approaches the set of equilibrium points $(\mathbf{c}^*, \mathbf{w}^*)$ as $t \rightarrow \infty$. \square

5.4 Conclusions

In this chapter, for the alternative optimisation problem (4.1), upper and lower bounds on the sensitivity of optimal path prices $\mathbf{q}^*(\mathbf{w})$ and aggregate source rates $\{x_s^*(\mathbf{w})\}$ to the variations of utility weight coefficients \mathbf{w} are derived. The sensi-

tivity results show that for each source $s \in S$, optimal path prices q_s^* grow as its weight w_s increases.

Given the correlation between path queueing delays and path prices $\mathbf{q}^*(\mathbf{c})$ in the scheduling algorithms (4.13)-(4.14), an alternative approach is then proposed where utility weight coefficients are used as control variables to regulate end-to-end queueing delays. Based on this approach, an integral controller is incorporated in the scheduling algorithms (4.13)-(4.14) whereby each source regulates the queueing delay on its paths at the desired level, using its weight coefficient as the control variable.

The proposed integral controller is distributed, since it is implemented at each source and uses only local path queueing delay information. Moreover, since the equilibrium of the scheduling algorithms (4.13)-(4.14) combined with integral controller (5.5) is the optimal solution of the optimisation problem (4.1) with equilibrium weight coefficients \mathbf{w}^* , it results in maximal link utilisation. For future work, linearisation and linear system design methods can further be used to adjust the controller parameter for the desired transient behaviour. Thus the proposed solution meets the objectives stated in Section 1.3.4.

Finally, the conditions under which algorithms (4.13)-(4.14) combined with the proposed integral controller (5.5) achieve asymptotic regulation of end-to-end delay are examined. The performance characteristics of the proposed solutions will also be demonstrated using simulation in the next chapter.

Chapter 6

Simulation Results

6.1 Introduction

In this chapter simulation experiments are performed to address three fundamental questions. Firstly, to illustrate that algorithms (4.13)-(4.14) converge to the optimal solutions of (4.8), despite using approximate values of link prices $\lambda(\mathbf{c})$ computed by the inner layer rate control algorithm (4.10). Secondly, to illustrate that algorithms (4.13)-(4.14) combined with (5.5) can regulate packet end-to-end latency, using an estimation of end-to-end delay as feedback in (4.13)-(4.14), and to compare their performance against the previously proposed main approaches to support delay-sensitive traffic. Finally, to assess the dynamic behaviour of the proposed algorithms when network configuration changes.

In Section 6.2, the simulated network and its mathematical model are described. The SimEvents implementation of the proposed algorithms in Chapters 4 and 5 for the network model is described in Section 6.3. The result of the simulation experiments is presented in Section 6.4. The conclusions are presented in Section 6.5.

6.2 Network Model

For simulation experiments the network topology in Figure 6.1 is considered, where there are two source-destination pairs $A \rightarrow C$ and $E \rightarrow D$. For source-destination pair $A \rightarrow C$, there are two alternative paths $A \rightarrow B \rightarrow C$ and $A \rightarrow D \rightarrow C$. Consequently, the routing matrix for the network is given by

$$R = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

Each active link is assumed to have a fixed data rate of c_0 packet per second. To model the scheduling constraint (2.3), the notions of contention graph and contention matrix [7] are used. In the contention graph, vertices represent links and edges represent the contention between the links. Maximal cliques of the contention graph embody the local contention among links; Links that belong to the same maximal clique cannot be active simultaneously. Let N be the number maximal cliques in the contention graph. The $N \times L$ contention matrix F is then defined by

$$F_{n,l} = \begin{cases} \frac{1}{c_0} & \text{if link } l \in L \text{ belongs to the maximal clique } n, \\ 0 & \text{otherwise.} \end{cases}$$

Thus, a necessary condition for scheduling is given by

$$F\mathbf{c} \leq \mathbf{1} \tag{6.1}$$

It can be shown that (6.1) is also a sufficient condition for scheduling if the contention graph is perfect [7]. Thus, in this example (6.1) models the scheduling constraint (2.3).

It is assumed that each wireless node can only communicate with one other node at any time. This results in the contention graph shown in Figure 6.2. There are four maximal cliques: links (3,4,5), links (1,2), links (1,3) and links (2,4).

Thus, only one link on each path of source-destination pair $A \rightarrow C$ can be active. Moreover, links on path $A \rightarrow D \rightarrow C$ are in contention with the link on the only path of source-destination pair $E \rightarrow D$, since they all share node D . The contention matrix for the network is then given by

$$F = \begin{bmatrix} 0 & 0 & 1 & 1 & 1 \\ 1 & 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 \end{bmatrix}$$

Since the contention graph in Figure 6.2 has no odd holes, it is perfect and therefore (6.1) is a sufficient scheduling constraint in this case.

The utility functions for both sources are assumed to be of the form $w_s f_s(x_s) = w_s \ln(x_s)$. Logarithmic utility functions have all the necessary properties described in Section 2.2, i.e. twice continuously differentiable, strictly concave, increasing, and have strictly positive second derivative. Moreover, they have been shown to achieve weighted proportionally fair resource allocation, i.e. any deviation from optimal rate allocation results in less than, or equal to zero, weighted sum of proportional changes to each source's rate [35].

6.3 Implementation Using SimEvents

As the main objective of the simulation is to evaluate packet end-to-end latency when the proposed algorithms in Chapters 4 and 5 are implemented in a multihop wireless network architecture, the discrete-event simulation software SimEvents is used for simulation experiments. The SimEvents simulation models are presented in the Appendix. Figure A.1 shows the simulation model of the network in Figure 6.1. Nodes buffers are modelled as FIFO queues with infinite capacities. Links are modelled as servers with packet service time computed by the scheduling subsystem. Packets on each path are generated by the entity generators with entity intergeneration time controlled by the Source Rate Control subsystems. Packet end-to-end latencies are recorded by timers on every path.

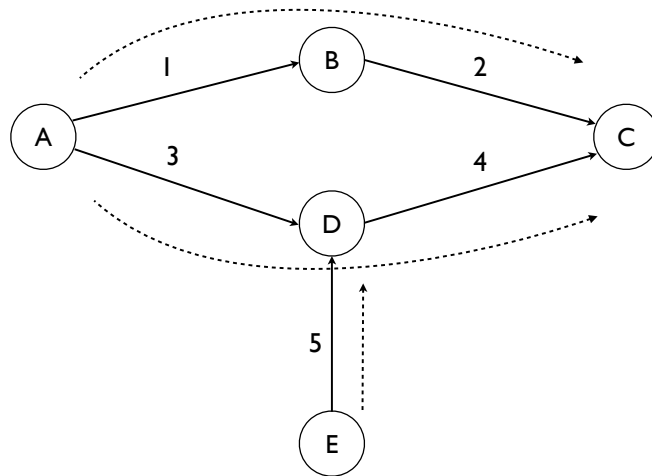


Figure 6.1: Network topology and alternative paths for source-destination pairs $A \rightarrow C$ and $E \rightarrow D$

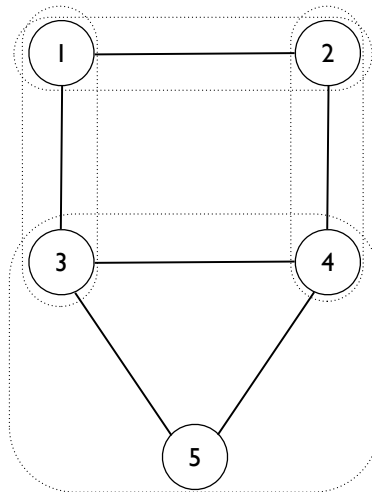


Figure 6.2: Network contention graph and its maximal cliques

The Scheduling subsystem implements algorithms (4.13)-(4.14), using link prices computed by Link Price Update subsystems. As indicated in Figure A.1, the Scheduling subsystem runs at discrete time steps with the duration of 500 time units, in order to allow link prices and path rates to reach near their optimal level for the current link data rates. To examine the robustness of the proposed algorithms in delay regulation in the presence of discrepancies between link prices and the corresponding average queueing delays, it is assumed that link prices are rather an estimation of β multiple of link average queueing delays, based on duality-based algorithm (4.10) initialised according to

$$\lambda_l^0(\mathbf{c}(k)) = \begin{cases} \frac{c_l(k-1)}{c_l(k)} \lambda_l(\mathbf{c}(k-1)) & k > 0 \\ 0 & k = 0 \end{cases} \quad (6.2)$$

where $\lambda_l^0(\mathbf{c})$ is the initial point of (4.10) given \mathbf{c} and k corresponds to the discrete time step of the scheduling algorithm. Link Price Update subsystems implement duality-based algorithm (4.10), using path rates computed by Source Rate Control subsystems. They are initialised according to (6.2) every time link data rates are updated by the Scheduling subsystem, and run in continuous-time for the duration of scheduling time step, as shown in Figure A.3. Source Rate Control subsystems compute path rates in (4.3), given fixed link data rates. However, since Source 1 has two alternative paths, the computation of path rates in the Source 1 Rate Control subsystem is based on inclusion of small quadratic term $\delta \mathbf{x}^T \mathbf{x}$ in (4.3), in order to stabilise the path rates (Figure A.2). Source Rate Control subsystems run jointly with Link Price Update subsystems in continuous-time for the duration of scheduling time step. Source Delay Regulator subsystems implement algorithm (5.5), using current average end-to-end delays as feedback. Source Delay Regulator subsystems run at discrete time steps with the duration of 500 time units, simultaneously with the Scheduling subsystem, in order to allow link prices and path rates to reach near their optimal level for the current link data rates.

The algorithms implemented in the simulation subsystems are then summarised as follows.

Scheduling subsystem At step k :

$$\mathbf{c}(k) = \mathbf{c}(k-1) + \gamma (\tilde{\mathbf{c}} - \mathbf{c}(k-1)) \quad (6.3)$$

where

$$\tilde{\mathbf{c}} = \arg \max_{\boldsymbol{\varsigma}} \boldsymbol{\lambda}(\mathbf{c}(k-1))^T \boldsymbol{\varsigma} \quad \text{s.t.} \quad F\boldsymbol{\varsigma} \leq \mathbf{1} \quad (6.4)$$

and $0 < \gamma \leq 1$.

Source Delay Regulator subsystem At step k :

$$w_s(k) = \left[w_s(k-1) + \alpha \left(d_s - \frac{q_s(\mathbf{c}(k-1), \mathbf{w}(k-1))}{\beta} \right) \right]^+ \quad \forall s \in S \quad (6.5)$$

where $\alpha > 0$, and $q_s(\mathbf{c}(k-1), \mathbf{w}(k-1))$ is computed using the values of link prices (computed by Link Price Update subsystem) at the end of the previous step $k-1$.

Link Price Update subsystem Given the current link data rates $\mathbf{c}(k)$, continue until the next scheduling update at step $k+1$:

$$\dot{\lambda}_l = \frac{\beta}{c_l(k)} [R_l \mathbf{x}(\boldsymbol{\lambda}(t)) - c_l(k)]_{\lambda_l}^+ \quad \forall l \in L \quad (6.6)$$

where $\beta > 0$, the initial point λ_l^0 is given by (6.2), and $\mathbf{x}(\boldsymbol{\lambda}(t))$ is computed by the Source Rate Control subsystem.

Source Rate Control subsystem Given $\boldsymbol{\lambda}(t)$ and $\mathbf{w}(k)$, solve the following systems of equations for all $s \in S$:

$$w_s(k) f'_s \left(\sum_{i \in I_s} x_i^s \right) - q_i^s(t) + 2\delta x_i^s = 0 \quad \forall i \in I_s \quad (6.7)$$

here $f_s(x_s) = \ln(x_s)$, and $\delta > 0$ is a small constant.

6.4 Results

The objective of the simulation experiments is threefold. Firstly, to illustrate that algorithms (6.3)-(6.4) indeed converge to the unique optimal solutions of (4.8), where link prices $\lambda(\mathbf{c}(k))$ are computed using (6.6)-(6.7) and (6.2) running for a finite time. Moreover, to observe how closely actual average end-to-end delays correspond to the path prices $q(\mathbf{c}(k))$ computed using (6.6)-(6.7) and (6.2). Secondly, to examine the robustness and accuracy of algorithms (6.3)-(6.4) combined with (6.5) in regulating packet end-to-end latency, where average end-to-end delay is used as feedback in (6.5) but link prices $\lambda(\mathbf{c}(k))$ in (6.3)-(6.4) are rather an estimation of β multiple of link average queueing delays based on (6.6)-(6.7) and (6.2). Furthermore, to compare their performance against the previously proposed main approaches to support delay-sensitive traffic, namely, the virtual data rates approach (Section 3.4) and the proposed approach by [24] (Section 3.5). Since the convergence of the proposed solutions were shown on basis of the assumption that network topology, channel conditions and other network configuration remain fixed, or their changes can be compensated at the physical layer, over the timescale of the problem (as discussed in Sections 1.2.6 and 1.3.2), the final objective is to assess the dynamic behaviour of the proposed algorithms in the presence of network configuration changes, specifically when a new flow enters the network.

In all experiments it is assumed that $c_0 = 1$ packet per milliseconds, $\gamma = 1 \times 10^{-2}$, $\beta = 1 \times 10^{-3}$, and $\delta = 1 \times 10^{-1}$. In the first experiment, source weights are fixed at $w_1 = 2$ and $w_2 = 1$ and algorithms (6.3)-(6.4) are simulated. The evolution of path transmission rates and link data rates are shown in Figures 6.3 and 6.4, respectively. Despite the fact that (6.3)-(6.4) use approximate values of link prices $\lambda(\mathbf{c}(k))$, link data rates converge to their optimal level (0.5, 0.5, 0.1667, 0.1667, 0.6667) relatively quickly, but have a non-smooth curve due to discontinuous nature of the right-hand side of original scheduling algorithm (4.13). The approximate values of link prices $\lambda(\mathbf{c}(k))$ result in slight oscillations in the path prices as shown in Figure 6.5, which in turn lead to the oscillations of the corresponding path rates (Figure 6.3). Although the rate oscillations for the only path of source 2 are also modest, they are more significant for both paths of

source 1. As described in Section 4.2, for sources with multiple alternative paths, only paths that have the minimum price will have positive flows at optimality. Thus, although path transmission rates converge quickly to a close neighbourhood of their optimal values (0.5, 0.1667, 0.6667), the almost coincidence of path prices for source 1 at non-differentiability points of link data rate curves results in equal path rates at the corresponding points. At the other points, however, due to the slight divergence of path prices, path rates for source 1 also diverge. In this case, the added quadratic term to (4.3) stops the rate of path 2 from becoming zero.

In Figure 6.5, path prices are compared with normalised end-to-end delays. Although the dynamic behaviour of path prices and the corresponding normalised end-to-end delays are similar, the normalised end-to-end delays lag slightly behind the path prices. Moreover, the magnitude of path prices for source 1 is approximately twice as high as their corresponding normalised end-to-end delays. The latter discrepancy is due to the fact that algorithms (6.6)-(6.7) leads to the same price, and consequently same data rate, for links belonging to the same path. In the simulation model, however, packet arrival rates at links 2 and 4 are equal to the data rates of links 1 and 3, respectively. Furthermore, since packet arrival rates and data rates of links 2 and 4 are almost equal throughout the simulation, their queueing delay is zero. Hence, the normalised end-to-end delay is almost half of the total price for each path of source 1.

In the second experiment, algorithms (6.3)-(6.4) are simulated jointly with (6.5) with delay bounds $d_1 = d_2 = 1 \times 10^3$ milliseconds, and $\alpha = 2 \times 10^{-5}$. Moreover, to compare their performance against the previously proposed main approaches to support delay-sensitive traffic, algorithms (6.3)-(6.4) are also simulated with the same setup as the first experiment but the dual-based algorithm (6.6) is modified using the following alternative methods

- virtual data rates (Section 3.4):

$$\dot{\lambda}_l = \frac{\beta}{c_l(k)} [R_l \mathbf{x}(\boldsymbol{\lambda}(t)) - \rho c_l(k)]_{\lambda_l}^+ \quad \forall l \in L \quad (6.8)$$

where $\rho = 0.99$ to retain high link utilisation.

- the proposed approach by [24] (Section 3.5):

$$\dot{\lambda}_l = \frac{\beta}{c_l(k)} \left[R_l \mathbf{x}(\boldsymbol{\lambda}(t)) - \theta_l^{-1} \left(\frac{\lambda_l(t)}{R_l \mathbf{b}} \right) \right]_{\lambda_l}^+ \quad \forall l \in L \quad (6.9)$$

where $\theta_l(c_l, y_l)$ is computed from (2.6) with $c_l = c_l(k)$, \mathbf{b} is an $I \times 1$ vector with elements $b_i^s = b_s$, $\forall i \in I_s$, and b_s indicates the degree of sensitivity of source s to delay. Here, $b_1 = 0.001$ and $b_2 = 0.0005$ to retain high link utilisation.

The evolution of packet end-to-end delays in both first and second experiments are compared in Figure 6.6. Despite the discrepancies between path prices and actual average end-to-end delay, such as time lag and difference in magnitude, using actual average end-to-end delay as feedback and with proper choice of parameter α , algorithms (6.3)-(6.4) combined with (6.5) regulate end-to-end delays at their upper bound levels with good precision. Moreover, since their equilibrium is the solution of the optimisation problem (4.1) with equilibrium parameter \mathbf{w}^* , they lead to maximal link utilisation. The delay regulator's parameter α can further be adjusted to achieve the desired transient behaviour of end-to-end delays by, for example, by linearisation around the equilibrium point and using linear control systems design methods.

Other alternative approaches reduce the delay to near zero in the long term, as they lead to under-utilised links at the equilibrium (Section 3.6). This result confirms the inaccuracy of the delay estimation using $M/D/1$ queue delay models discussed in Section 2.4. Moreover, the rate by which they reduce delay is inversely correlated to the levels of link utilisation at their equilibrium. Specifically, in the approach based on virtual link data rates, the speed of delay reduction increases as parameter ρ , and hence link utilisation, decreases. Similarly, in the approach proposed by [24], the rate of delay reduction increases as the delay-sensitivity coefficients b_s , $s \in S$, increase. Consequently, the delay functions, which grow exponentially as link utilisations approach unity, have higher weight in the objective function (Section 3.5), and as a result link utilisation at optimality further decreases. In addition, the other alternative approaches are incapable of regulating

end-to-end delays at pre-specified levels with desired transient behaviour. Nevertheless, these examples show that the proposed scheduling algorithms (4.13)-(4.14) can be effectively used in conjunction with other approaches designed to support delay sensitive traffic, in particular the rate control approach proposed by [24] which is originally designed for networks supporting heterogeneous traffic where link data rates are fixed and sources transmit only one flow using a fixed path.

The final experiment simulates the response of algorithms (6.3)-(6.4) combined with (6.5) when a new flow enters the network. Specifically, the second experiment is repeated with the difference that flow $E \rightarrow D$ starts at simulation time 2.5×10^4 . The evolution of path transmission rates, link data rates, path prices, packet end-to-end delays and source weights are shown in Figures 6.7-6.11. Before the start of flow $E \rightarrow D$, path rates of flow $A \rightarrow C$ as well as data rates of links 1-4 have approached near their maximum levels (0.5,0.5) packets/milliseconds. After the start of flow $E \rightarrow D$, the transmission rate of path 1 of flow $A \rightarrow C$ and the data rates of associated links continue to approach their maximum level 0.5 packets/milliseconds. However, as the transmission rate of flow $E \rightarrow D$ rises, the data rate of link 5 increases causing the data rates of links 3 and 4, and hence the transmission rate of path 2 of flow $A \rightarrow C$, to decrease. At the start of flow $E \rightarrow D$, its transmission rate shoots up rapidly, resulting in a big surge in its packet end-to-end delay. This causes its weights to go down rapidly, which in turn leads to a quick drop in its packet end-to-end delay. Consequently, packet end-to-end delays approach their upper bounds quickly. After this point, there is a relatively long period of slight oscillations before stabilisation, due to the feedback errors discussed earlier. The oscillations of path rates and link data rates have higher magnitude, due to the extra disturbances caused by the perturbations in parameters \mathbf{w} . This experiment shows that in the presence of network configuration changes, the proposed algorithms converge only if they converge at a shorter timescale than the changes in network configuration.

6.5 Conclusions

In this chapter an example of simultaneous transmissions of several delay-sensitive traffics over a multi-hop wireless network is formulated as the original optimisation problem (2.1)-(2.5). The example network is modelled as a queueing network where link data rates and source data transmission rates are controlled by the proposed algorithms in Chapters 4 and 5, and subsequently implemented using SimEvents discrete-event simulation software.

The first experiment shows that despite using approximate values of link prices $\lambda(\mathbf{c}(k))$, scheduling algorithms (6.3)-(6.4) converge relatively quickly to the unique optimal solutions of (4.8), but the evolution of link data rates is non-smooth due to discontinuous nature of the right-hand side of the original scheduling algorithm (4.13). The approximate values of link prices $\lambda(\mathbf{c}(k))$ result in slight oscillations in the path prices, which in turn lead to the oscillations of the corresponding path rates. The resulting path rate oscillations are greater for sources with multiple paths due to the well-known oscillation problem of rate control algorithms (6.6)-(6.7).

The second experiment shows that despite the presence of feedback error caused by the delay approximation used in the rate control algorithms (6.6)-(6.7), algorithms (6.3)-(6.4) combined with (6.5) regulate end-to-end delays at their upper bound levels with good precision. They also enable the control of transient behaviour of end-to-end delays using standard control systems design techniques. Moreover, since their equilibrium is the solution of the optimisation problem (4.1) with equilibrium parameter \mathbf{w}^* , by design they lead to maximal link utilisation. On the other hand, alternative approaches such as virtual data rates (Section 3.4), and the proposed approach by [24] (Section 3.5), can only achieve relatively high (but still not maximal) link utilisation at expense of slow reduction of end-to-end delays. In addition, they cannot regulate end-to-end delays at pre-specified levels with desired transient behaviour. In fact they lead to near zero delays in the long term, confirming the inaccuracy of the delay estimation using $M/D/1$ queue delay models. This experiment, however, demonstrates that the proposed scheduling algorithms (6.3)-(6.4) can be effectively used in conjunction with other approaches

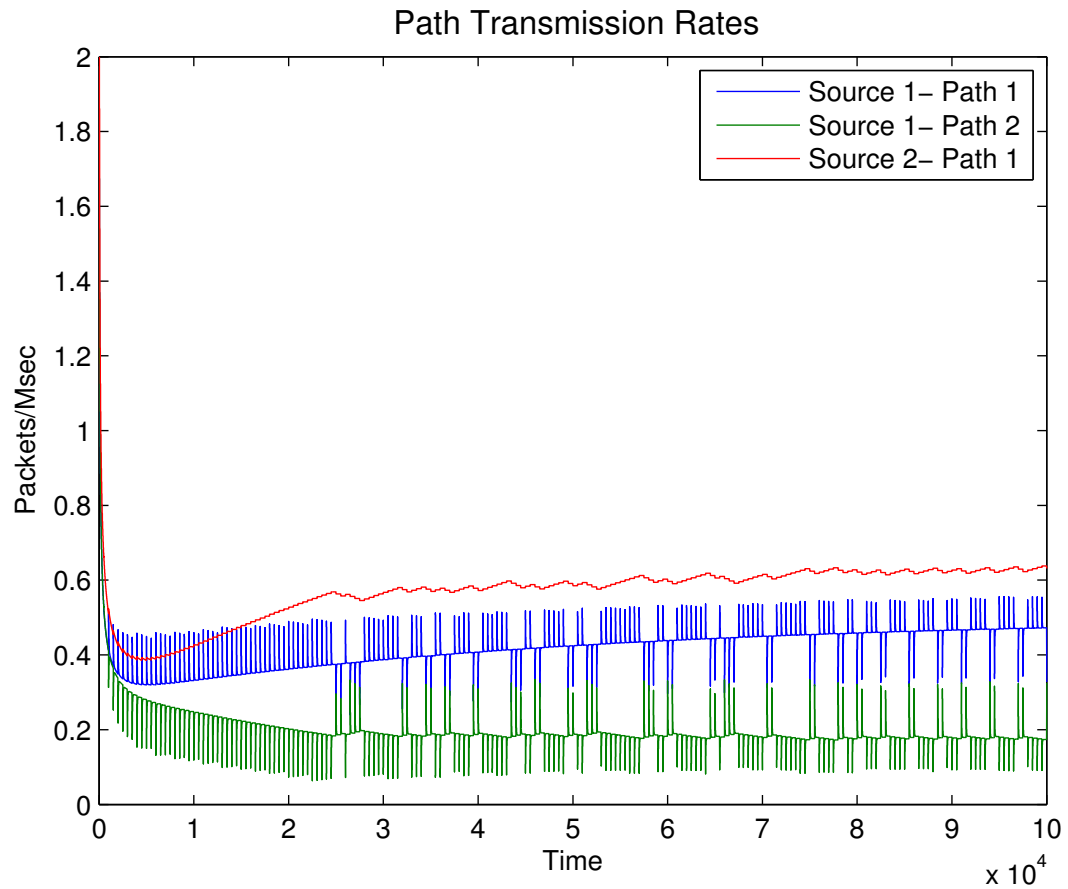


Figure 6.3: Path transmission rates when source weights are fixed at $w_1 = 2$ and $w_2 = 1$ and algorithms (4.13)-(4.14) are simulated

designed to support delay sensitive traffic, in order to reduce delay.

The third experiment demonstrates the limitations of the proposed algorithms in the presence of network configuration changes, in that the proposed algorithms converge only if their convergence rate is faster than the rate of changes in network configuration.

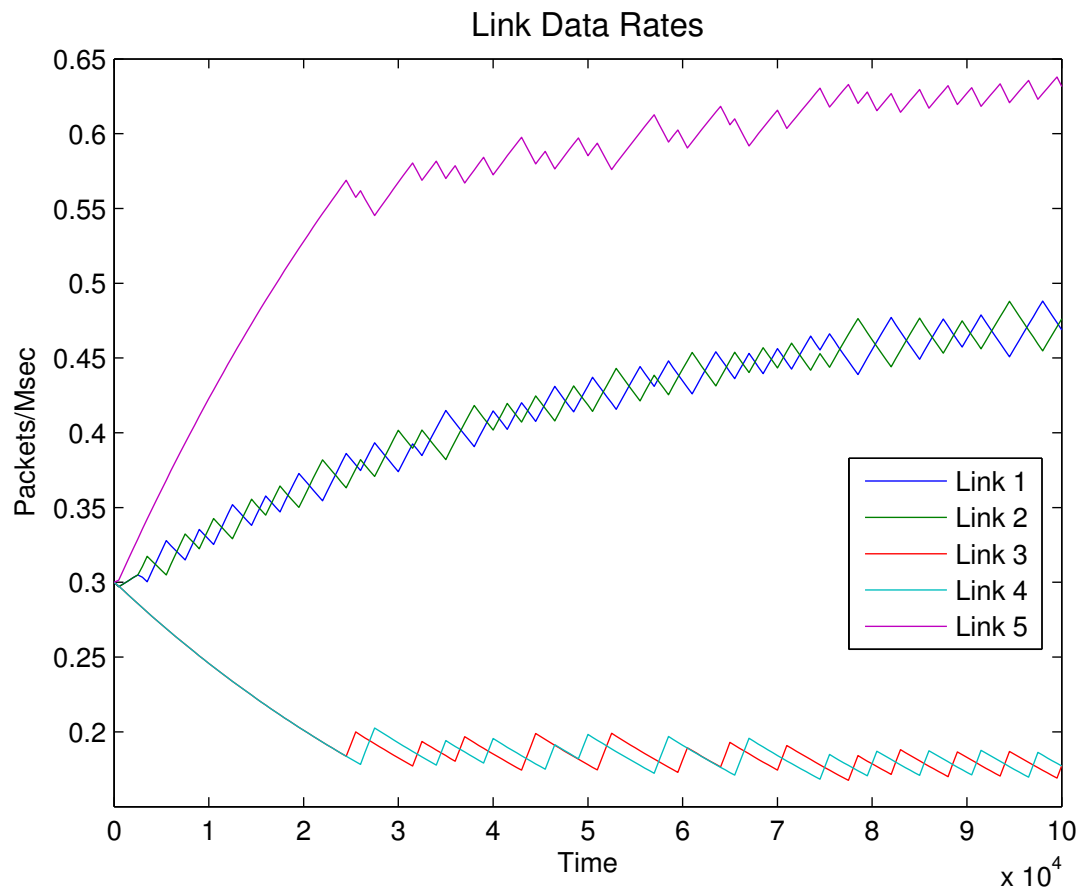


Figure 6.4: Link data rates when source weights are fixed at $w_1 = 2$ and $w_2 = 1$ and algorithms (4.13)-(4.14) are simulated

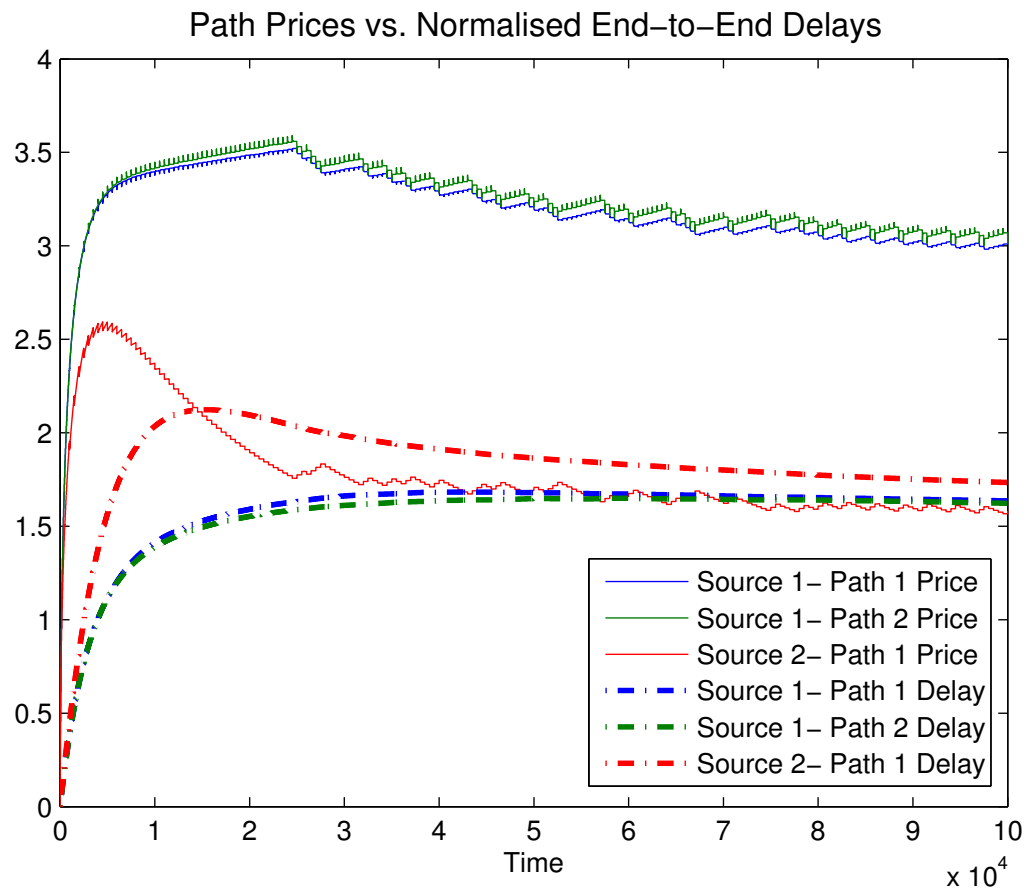


Figure 6.5: Path prices and normalised end-to-end delays when source weights are fixed at $w_1 = 2$ and $w_2 = 1$ and algorithms (4.13)-(4.14) are simulated

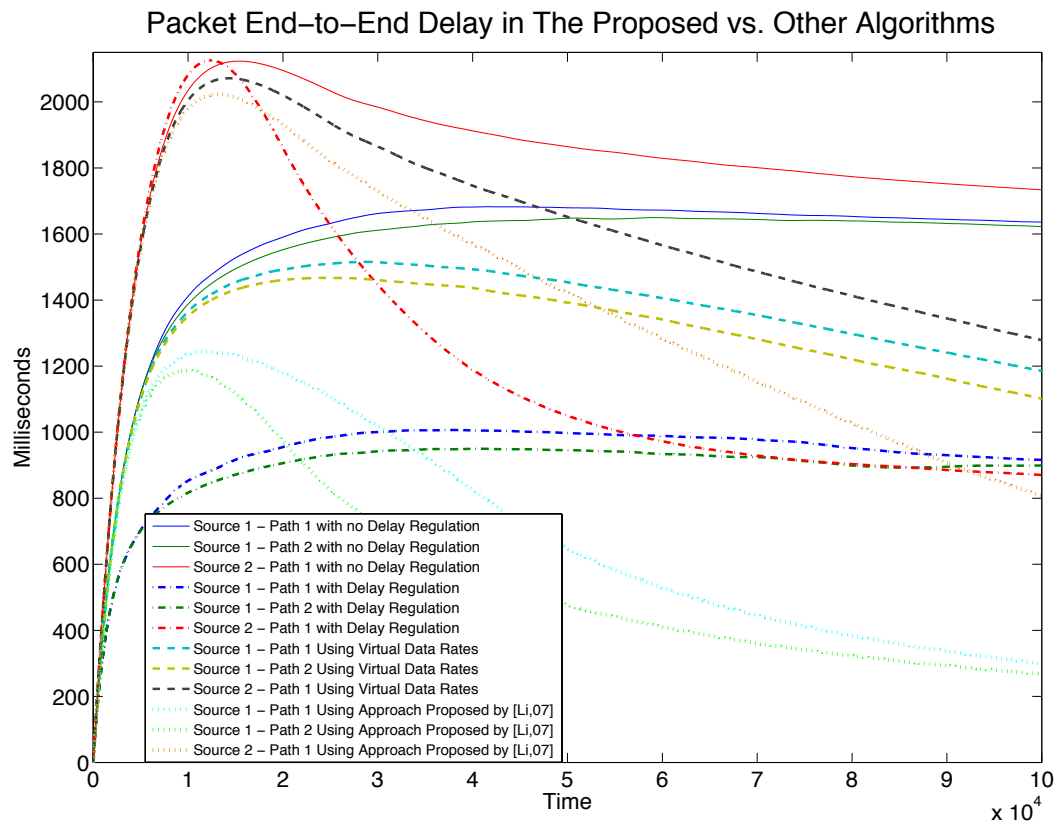


Figure 6.6: Packet end-to-end delay when algorithms (4.13)-(4.14) are simulated jointly with (5.5) with delay bounds $d_1 = d_2 = 1 \times 10^3$ msec

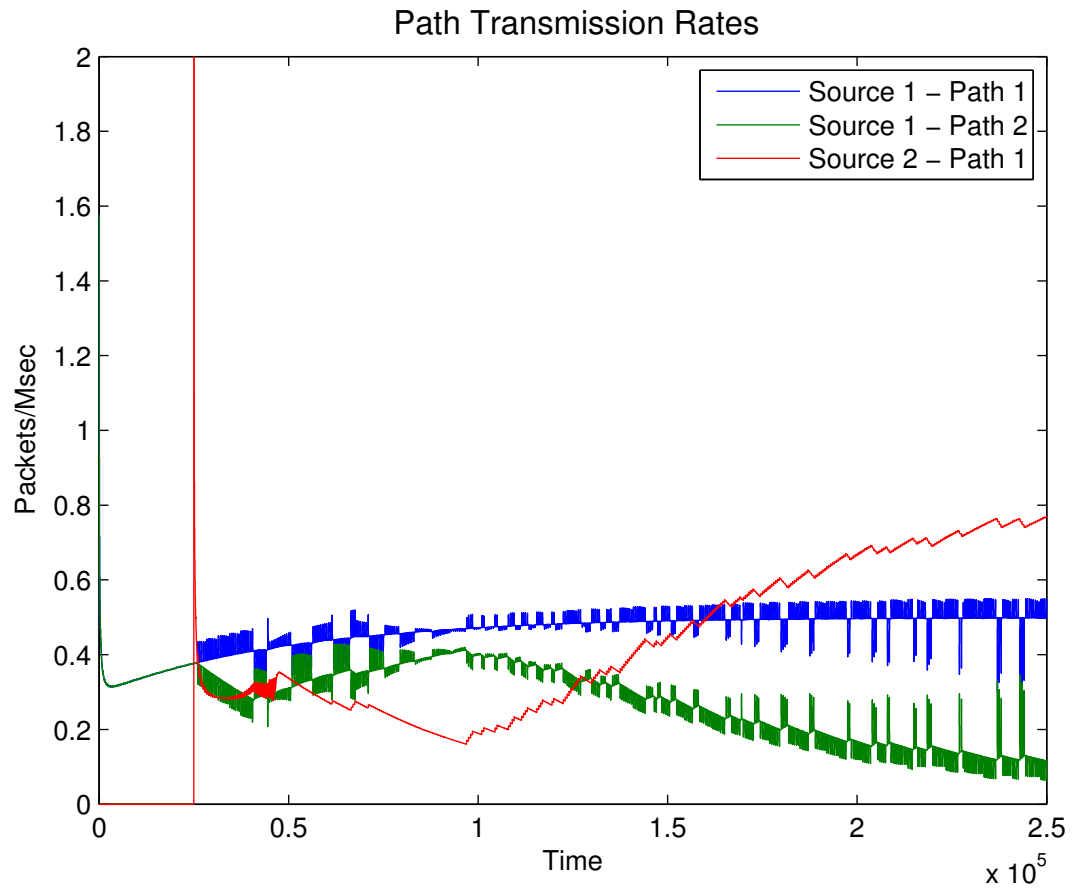


Figure 6.7: Path transmission rates when (4.13)-(4.14) are simulated jointly with (5.5) with delay bounds $d_1 = d_2 = 1 \times 10^3$ msec, and flow $E \rightarrow D$ starts at time 2.5×10^4

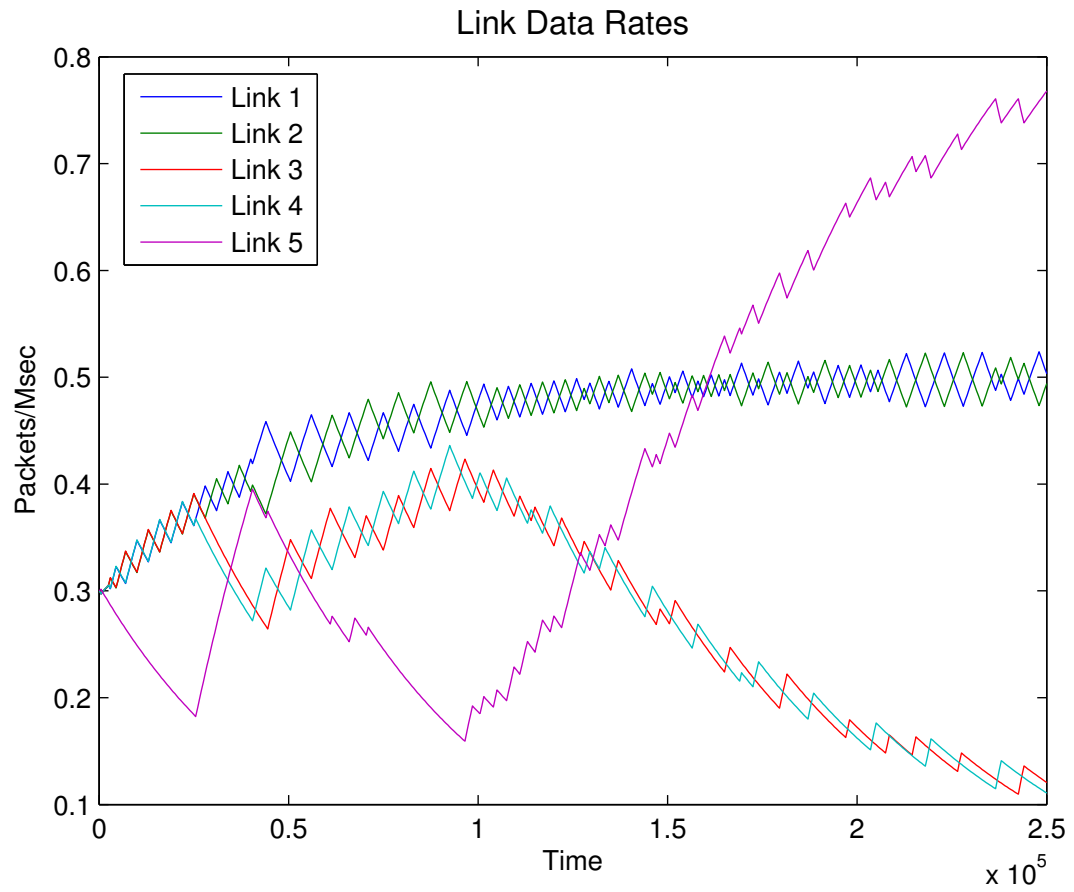


Figure 6.8: Link data rates when (4.13)-(4.14) are simulated jointly with (5.5) with delay bounds $d_1 = d_2 = 1 \times 10^3$ msec, and flow $E \rightarrow D$ starts at time 2.5×10^4

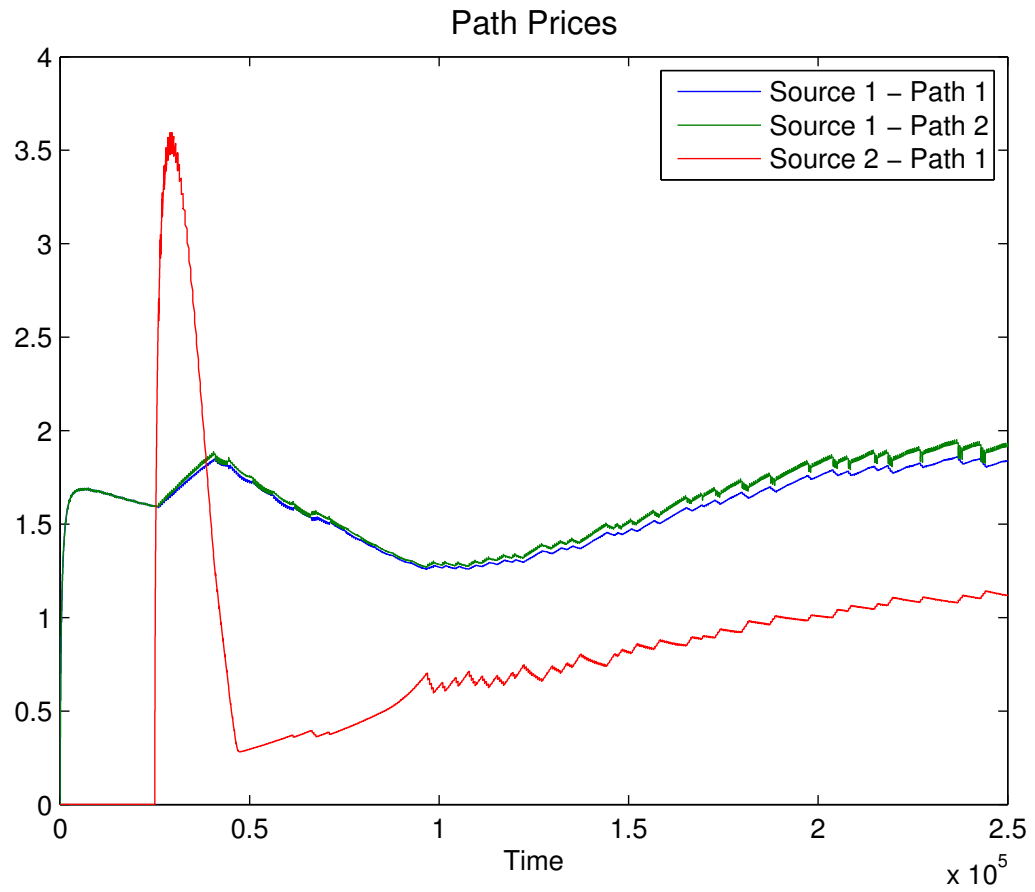


Figure 6.9: Path prices when (4.13)-(4.14) are simulated jointly with (5.5) with delay bounds $d_1 = d_2 = 1 \times 10^3$ msec, and flow $E \rightarrow D$ starts at time 2.5×10^4

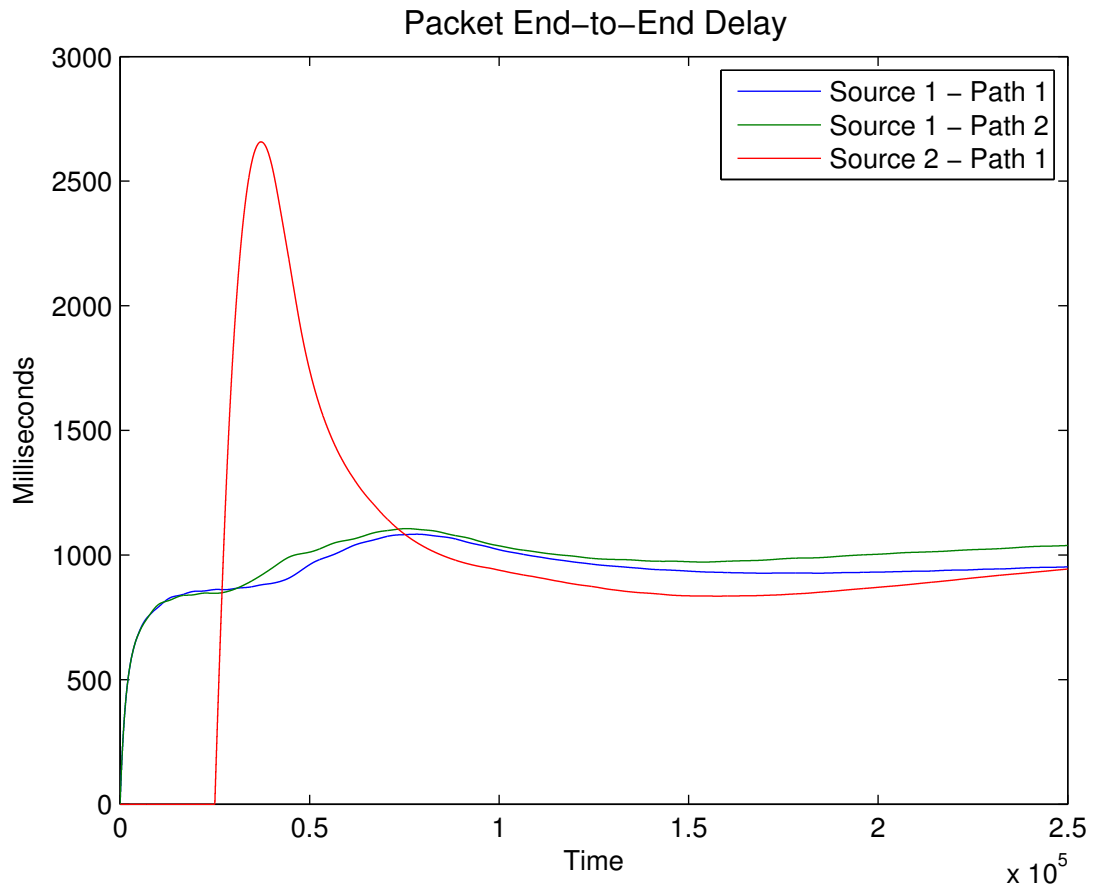


Figure 6.10: Packet end-to-end delays when (4.13)-(4.14) are simulated jointly with (5.5) with delay bounds $d_1 = d_2 = 1 \times 10^3$ msec, and flow $E \rightarrow D$ starts at time 2.5×10^4

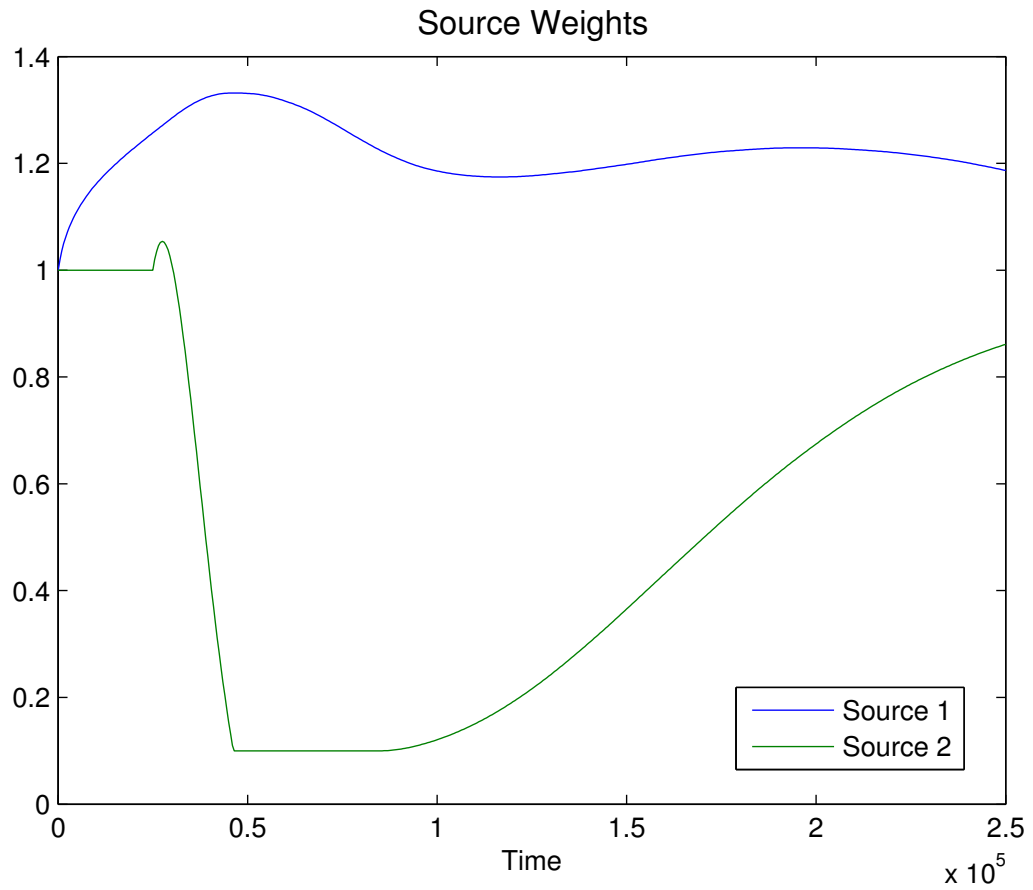


Figure 6.11: Source weights when (4.13)-(4.14) are simulated jointly with (5.5) with delay bounds $d_1 = d_2 = 1 \times 10^3$ msec, and flow $E \rightarrow D$ starts at time 2.5×10^4

Chapter 7

Conclusions

7.1 Main Findings

The main focus of this research is the problem of supporting delay-sensitive traffic with elastic data rate requirements and hard end-to-end delay constraints in multi-hop wireless networks, using source data transmission rates and link data rates as the key design variables. It is assumed that routing tables are already computed by a source-driven routing algorithm, and remain constant within the time horizon of the problem. Furthermore, the set of feasible link data rates, or schedules, are known and remain constant over the time horizon of the problem. The network utility maximisation (NUM) framework is adopted as the main design method as it enables the design of distributed and efficient solutions.

The conventional modelling of the delay constraints based on $M/D/1$ queue approximation of links is shown to have several major flaws. Firstly, the key assumption behind $M/D/1$ queue delay model, i.e. Poisson arrival of packets at links, is invalid in this problem, and delay is mainly determined by the transient behaviour of the rate control and scheduling algorithms. Thus this model leads to inaccurate delay estimation. Secondly, it leads to inefficient utilisation of links at optimality since their estimated delay grows exponentially as link flow rates approach their capacities.

The joint rate control and scheduling NUM problem for the elastic traffic,

whose QoS is modelled as concave utility function of its transmission rate, has been shown to lead to the canonical distributed rate control algorithm, as well as scheduling algorithms, which in many cases, are simple, efficient and distributed. Modelling the QoS of delay-sensitive traffic as a non-concave utility function of its transmission rate, has also been shown to lead to the solutions with similar properties in many cases; however, the utility function in this approach only characterises QoS in terms of packet inter arrival delay rather than end-to-end queueing delay.

Moreover, previous work attempts to address end-to-end queueing delay requirements of delay-sensitive traffic have been mainly based on either reducing link utilisation, such as using virtual data rates and minimising network congestion, or approximation of links as $M/D/1$ queues. Both approaches result in unpredictable transient behaviour of packet delays, and inefficient link utilisation under optimal resource allocation.

The efficient and distributed nature of the solution algorithms, as well as the efficiency of the optimal solutions of the the joint rate control and scheduling NUM problem for the elastic traffic motivates an approach, in which in place of hard delay constraints based on inaccurate $M/D/1$ delay estimates, traffic end-to-end delay needs are guaranteed by some forms of concave and increasing utility functions of traffic source rates, similar to the utility functions of the elastic traffic.

An alternative NUM formulation is then considered where the delay constraints are omitted and the original concave utility functions are multiplied by weight coefficients. The alternative NUM problem is then transformed into a master scheduling problem and the well-known multi-path rate control with fixed link data rates subproblem. At the inner layer, a duality-based rate control algorithm using link queueing delays as feedback ensures that optimal dual variables, or link prices, stay in proportion to the link queueing delays. Conditions on the number of disjoint paths are derived that guarantee unique optimal path rates. The theoretical basis of this result can be further used in future work to design a multi-path rate control algorithm that avoids the well-known path rates oscillation problem associated with the duality-based algorithm, given the conditions on the number of disjoint paths are satisfied by the topology control algorithm. A distributed

algorithm for the master scheduling problem is then proposed, and is shown to converge to the optimal data rates. The proposed algorithm at its core solves a well-known scheduling problem, for which efficient and distributed solutions have been developed in several cases.

Having derived upper and lower bounds on the sensitivity of path prices to the variations of utility weight coefficients for each traffic source, it becomes apparent that optimal path prices for each source increase with its utility weight coefficient. Given the correlation between path queueing delays and path prices in the proposed scheduling algorithm, this reaffirms the idea of satisfying traffic end-to-end delay constraints using some form of concave utility function. The form of the concave utility function that ensures the end-to-end delay requirements in this approach, however, will be dependent on the actual network configuration. As such, an integral controller is incorporated in the scheduling algorithm whereby each source regulates the queueing delay on its paths at the desired level, using its weight coefficient as the control variable. Upper bound on the step size of the proposed integral controller is then derived that ensures the proposed joint scheduling algorithm and delay regulator achieve asymptotic regulation of end-to-end delay.

Simulation experiments show that the presence of feedback error, i.e. approximate values of path prices, the proposed scheduling algorithm converges to the optimal solution of the alternative optimisation. However, approximate values of path prices aggravates the well-known path rates oscillation problem associated with the duality-based algorithm. Moreover, despite the additional feedback error caused by the delay approximation in the rate control algorithms, the proposed joint scheduling algorithm and the integral controller regulate end-to-end delays at their desired levels with good precision. However, alternative approaches based on virtual data rates and $M/D/1$ delay estimates, can only achieve relatively high (but still not maximal) link utilisation at expense of slow reduction of end-to-end delays. They appear to have asymptotic zero delay, confirming the inaccuracy of $M/D/1$ delay estimates, with unpredictable transient behaviour. The experiments also demonstrate that the proposed scheduling algorithms can be effectively used in conjunction with other approaches designed to support delay sensitive traffic, in

order to reduce delay. Finally, simulation experiments illustrate that the proposed algorithms converge only if their convergence rate is faster than the rate of changes in network configuration.

7.2 Contributions to Knowledge

This thesis addresses the problem of supporting delay-sensitive traffic with elastic data rate requirements and hard end-to-end delay constraints in multi-hop wireless networks, with source data transmission rates and link data rates as the main design freedom. Previous network utility maximisation based approaches to tackle end-to-end queueing delay requirements of delay-sensitive traffic have been mainly based on either reducing link utilisation, or approximations using $M/D/1$ queue delay estimates. Both approaches suffer from unpredictable transient behaviour of packet delays, and inefficient link utilisation under optimal resource allocation.

Motivated by the simple, efficient and distributed nature of the solution algorithms, as well as the efficiency of the optimal solutions of the similar NUM problem for the elastic traffic, an alternative approach is proposed, in which in place of hard delay constraints based on inaccurate $M/D/1$ delay estimates, traffic end-to-end delay needs are guaranteed by proper forms of concave and increasing utility functions of traffic source rates, similar to the utility functions of the elastic traffic. The proposed approach is realised by a scheduling algorithm that runs jointly with an integral controller whereby each source regulates the queueing delay on its paths at the desired level, using its utility weight coefficient as the control variable. The proposed scheduling algorithm at each step solves a familiar scheduling problem, for which efficient and distributed solutions have been developed in several cases, and the well-known multi-path rate control problem, which is solved using distributed duality-based algorithms.

The proposed algorithms are simple and distributed, as they are, for most part, based on canonical distributed rate control and scheduling algorithms. Furthermore, as they are based on solving the alternative concave optimisation problem, they lead to maximal link utilisation. The proposed algorithms are shown, us-

ing both theoretical analysis and simulation, to achieve asymptotic regulation of end-to-end delay given the step size of the proposed integral controller does not exceed a specified upper limit. The step size of the proposed integral controller can potentially be further adjusted to achieve the desired transient behaviour, using linearisation and linear system design methods.

It is well known that for traffic sources with multiple alternative paths, path rates in duality-based algorithms do not converge and continuously oscillate. In this research conditions on the number of disjoint paths are derived that guarantee unique optimal path rates. The theoretical basis of this result can potentially be further used to design a duality-based multi-path rate control algorithm that converges the unique optimal path rates.

7.3 Limitations of the Work

There are several limitations associated with the proposed solution, which also provide motivations for possible future research directions.

Firstly, a key assumption in the proposed solution is that routing information at each traffic source, which are already computed by a source-driven routing, remain constant in the time horizon of the algorithm. This implies that variations at the physical layer can be compensated at the physical layer without affecting the network topology, and hence the routing information in the time horizon of the solution. Likewise, it is assumed that the set of feasible link data rates, or schedules, remain constant over the time horizon of the problem, which means that variations in link SINR levels can be compensated quickly by power control or adjusting other physical layer parameters without affecting the link data rates within the time horizon of the solution. While these are widespread assumptions in most related literature, they may not be valid in highly dynamic scenarios where network configuration changes rapidly.

Secondly, the proposed algorithm consists of an outer-layer scheduling algorithm and delay regulator, where each step of the outer layer algorithms requires the convergence of the inner-layer multi-path rate control algorithm. In order

for the outer-layer joint scheduling algorithm and delay regulator to converge, inner-layer multi-path rate control algorithm has to converge sufficiently close the optimal solutions. In addition to the slow convergence, distributed implementation of this two-layer structured algorithm is rather difficult, since it is difficult to control the accuracy of the solution of the inner-layer algorithm in a distributed setting.

7.4 Future Work

While the scope of this thesis was limited to the development and testing an alternative approach for regulating end-to-end queueing delay with greater efficiency compared to the previous NUM approaches, there remain several interesting directions for future research.

Firstly, given the dependency of the proposed solution on the assumptions that the network topology and the feasible link data rates remain unchanged over the time horizon of the solution, in other words, the solution algorithm convergence rate is faster than the rate of changes in network conditions, as described in the previous section, it would be interesting to analyse the stability of the solution in the presence of frequent variations in channel conditions and network configuration modelled by random processes.

Secondly, given the distributed implementation issues caused by the two-layer convergence structure of the proposed algorithm, as discussed in the previous section, further work would be needed on the practical distributed implementation of the proposed algorithm, particularly in analysing and controlling the accuracy of the solution.

Thirdly, as pointed out in Chapter 4, it is widely known that while the duality-based rate control algorithms always converge to a dual optimal solution, for sources with multiple alternative paths, path rates do not converge and continuously oscillate. In Chapter 4, the second-order sufficient conditions for a unique local maximising point of the multi-path rate control subproblem were subsequently used to derive conditions on the number of disjoint paths that guarantee unique

optimal path rates. This suggests a promising approach to address the path rate oscillation problem, by using the second-order sufficient conditions to design a multi-path rate control algorithm that converges the unique optimal path rates, in conjunction with a topology control algorithm that guarantees the required conditions on the number of disjoint paths.

Fourthly, designing efficient feedback mechanisms for the proposed algorithms as well as analysing the solution stability in the presence of the induced feedback delay remains an open issue to investigate.

Lastly, the application of linearisation and linear system design methods to adjust the delay regulator parameter, in order to achieve a desired transient behaviour of packet end-to-end delay, would be an interesting future step.

Appendix A

SimEvent Simulation Models

Sample Times for 'RateControlSchedulingDelayReg6aSim42Refined'

Annotation	Description	Value
Cont	Continuous	0
FiM	Fixed in Minor Step	[0,1]
D1	Discrete 1	1
D2	Discrete 2	500
V1	Variable 1	Unknown
H	Hybrid	Not Applicable

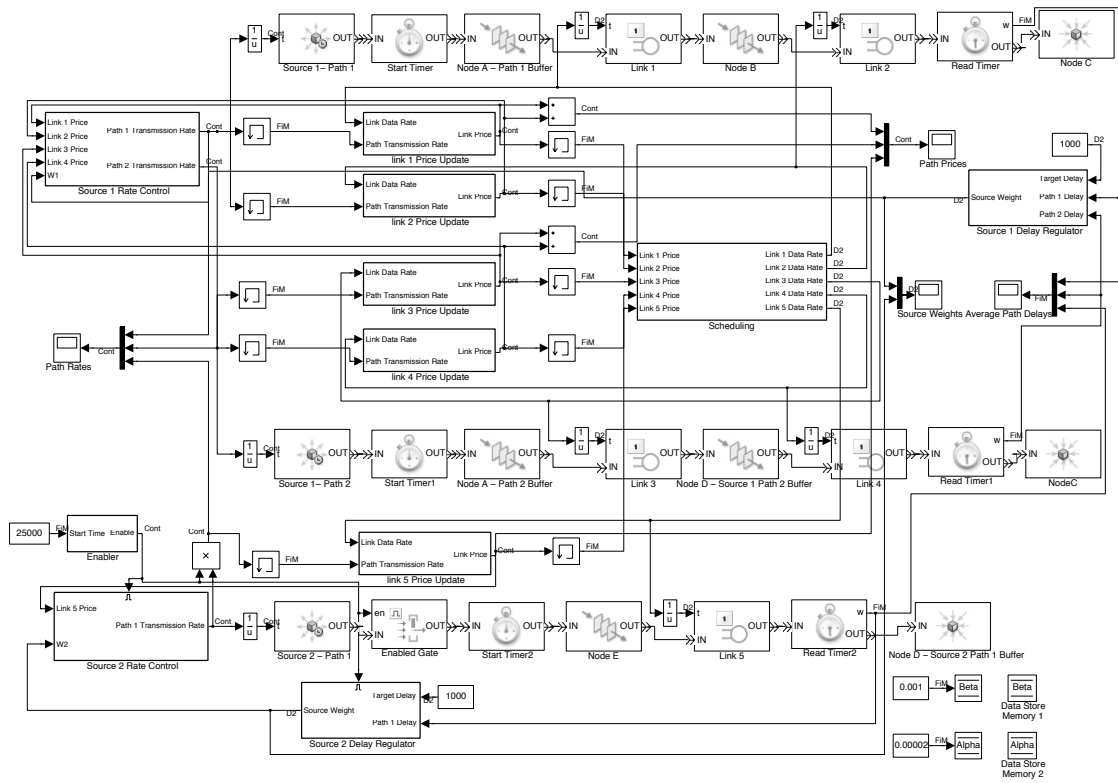


Figure A.1: Simulation model of the network in Figure 6.1

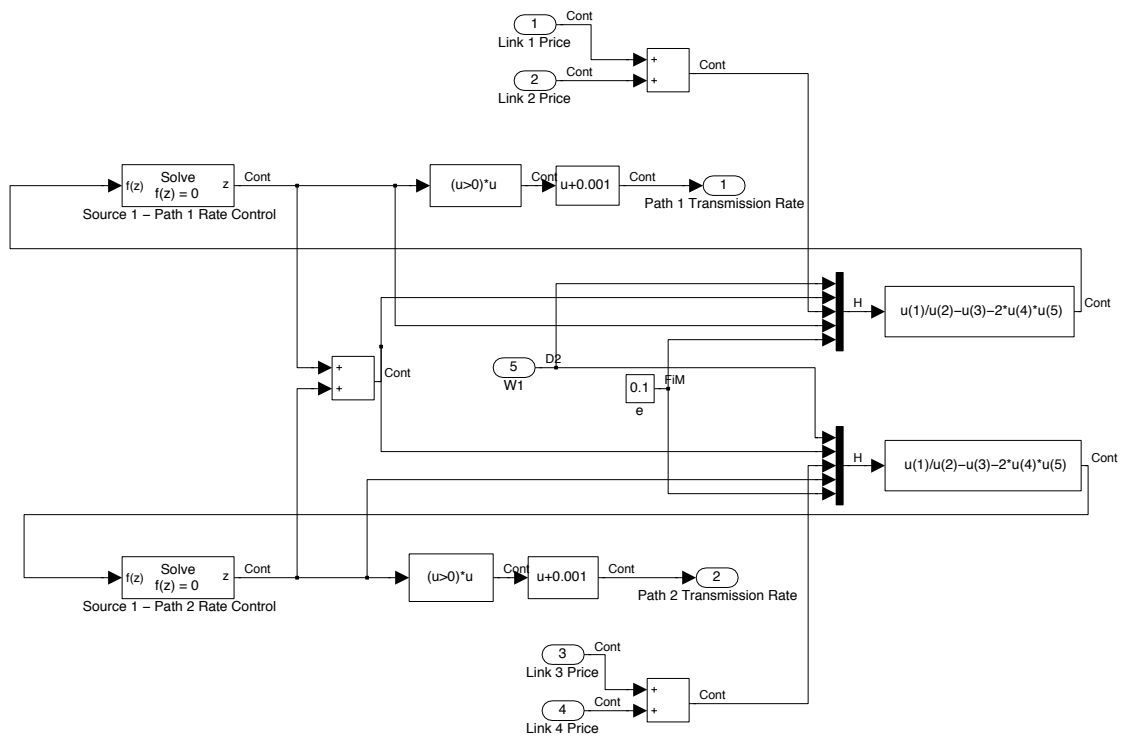


Figure A.2: Source 1 Rate Control subsystem

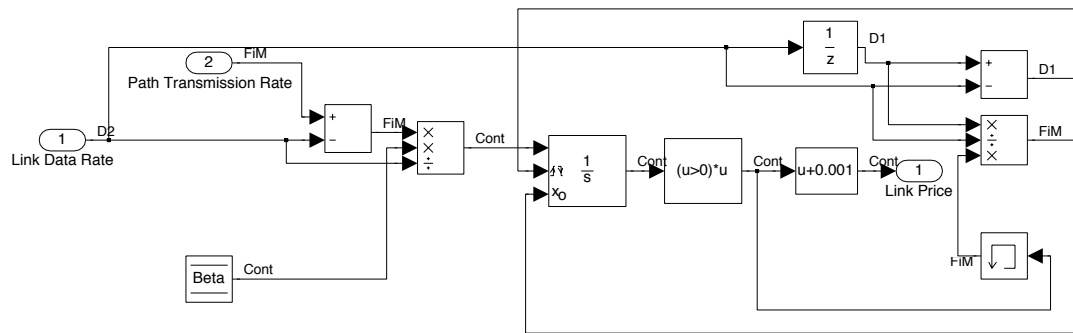


Figure A.3: Link Price Update subsystem

Appendix B

Mathematical Background

B.1 Sensitivity Analysis in Nonlinear Programming

Consider the following mathematical program

$$\begin{aligned} \min_x \quad & f(x, \epsilon) \\ \text{s.t.} \quad & x \in R(\epsilon) \\ & \epsilon \in T \end{aligned}$$

where $R(\epsilon)$ represents a constraint set as a function of the parameter ϵ , and X and T are topological vector spaces. Let $f^*(\epsilon) = \inf_x \{f(x, \epsilon) \mid x \in R(\epsilon)\}$, and define the mapping $S : T \rightarrow X$ by $S(\epsilon) = \{x \in R(\epsilon) \mid f(x, \epsilon) = f^*(\epsilon)\}$.

Theorem 2.2.6 [14] Assume that the space X is locally convex and that $R(\epsilon) \neq \emptyset$ for each $\epsilon \in T$ and that R is continuous and convex-valued on T (i.e., $R(\epsilon)$ is convex for each $\epsilon \in T$). If $f(x, \epsilon) = \min \{f_1(x, \epsilon), f_2(\epsilon)\}$, where f_1 is continuous on $X \times T$ and strictly quasi-convex in x for each fixed ϵ , and f_2 is continuous on T , then S is continuous and convex-valued on T .

Consider the problem of determining a local solution $x(\epsilon)$ of

$$\begin{aligned} \min_x \quad & f(x, \epsilon) \\ \text{s.t.} \quad & g_i(x, \epsilon) \geq 0, \quad i = 1, \dots, m \\ & h_j(x, \epsilon) = 0, \quad j = 1, \dots, p \end{aligned} \tag{B.1}$$

where $x \in \mathbf{R}^n$ and ϵ is a parameter vector in \mathbf{R}^k .

The Lagrangian of (B.1) is defined as

$$L(x, u, w, \epsilon) \equiv f(x, \epsilon) - \sum_{i=1}^m u_i g_i(x, \epsilon) + \sum_{j=1}^p w_j h_j(x, \epsilon) \tag{B.2}$$

Lemma 3.2.1 [14] (Second-order sufficient conditions for a strict local minimising point of problem (B.1).) If the functions defining problem (B.1) are twice continuously differentiable in a neighbourhood of x^* , then x^* is a strict local minimising point of problem (B.1) (i.e., there is a neighbourhood of x^* such that there does not exist any feasible $x \neq x^*$ where $f(x, 0) \leq f(x^*, 0)$) if there exist Lagrange multiplier vectors $u^* \in \mathbf{R}^m$ and $w^* \in \mathbf{R}^p$ such that the KKT conditions hold, i.e.

$$\begin{aligned} g_i(x^*, 0) &\geq 0, \quad i = 1, \dots, m \\ h_j(x^*, 0) &= 0, \quad j = 1, \dots, p \\ u_i^* g_i(x^*, 0) &= 0, \quad i = 1, \dots, m \\ u_i^* &\geq 0, \quad i = 1, \dots, m \end{aligned}$$

$$\begin{aligned} \nabla L(x^*, u^*, w^*, 0) &\equiv \\ \nabla f(x^*, 0) - \sum_{i=1}^m u_i^* \nabla g_i(x^*, 0) + \sum_{j=1}^p w_j^* \nabla h_j(x^*, 0) &= 0 \end{aligned}$$

and, further, if

$$\begin{aligned} z^T \nabla^2 L(x^*, u^*, w^*, 0) z &> 0 \quad \text{for all } z \neq 0 \text{ such that} \\ \nabla g_i(x^*, 0) z &\geq 0 \quad \text{for all } i, \text{ where } g_i(x^*, 0) = 0 \\ \nabla g_i(x^*, 0) z &= 0 \quad \text{for all } i, \text{ where } u_i^* > 0 \\ \nabla h_j(x^*, 0) z &= 0 \quad j = 1, \dots, p \end{aligned}$$

Remarks on Lemma 3.2.1 [14] It is shown that the conclusion of Lemma 3.2.1 can be strengthened to assert that x^* is a locally unique and hence isolated local minimum of problem (B.1), if the second-order sufficient conditions of Lemma 3.2.1 are strengthened by assuming these hold for all optimal Lagrange multipliers (u, w) associated with x^* and, given (B.1) is convex, the Slater's condition holds at x^* .

Theorem 3.2.2 [14] (First-order sensitivity results for a second-order local minimising point x^* .) If

1. the functions defining (B.1) are twice continuously differentiable in x and if their gradients with respect to x and the constraints are once continuously differentiable in ϵ in a neighbourhood of $(x^*, 0)$,
2. the second-order sufficient conditions for a local minimum of (B.1) hold at x^* , with associated Lagrange multipliers u^* and w^* ,
3. the gradients $\nabla g_i(x^*, 0)$ (for i such that $g_i(x^*, 0) = 0$) and $\nabla h_j(x^*, 0)$ (all j) are linearly independent,
4. $u_i^* > 0$ when $g_i(x^*, 0) = 0$, $i = 1, \dots, m$ (i.e. strict complementary slackness holds),

then

1. x^* is a local *isolated* minimising point of problem (B.1) and the associated Lagrange multipliers u^* and w^* are unique,
2. for ϵ in a neighbourhood of 0, there exists a unique, once continuously differentiable vector function $y(\epsilon) = (x(\epsilon), u(\epsilon), w(\epsilon))^T$ satisfying the second-order sufficient conditions for a local minimum of problem (B.1) such that $y(0) = (x^*, u^*, w^*)^T$, and hence $x(\epsilon)$ is a *locally unique* local minimum of problem (B.1) with associated unique Lagrange multipliers $u(\epsilon)$ and $w(\epsilon)$, and
3. for ϵ near 0, the set of binding inequalities is unchanged, strict complementary slackness holds, and the binding constraint gradients are linearly independent at $x(\epsilon)$.

B.2 Discontinuous Control

Consider the autonomous differential equation

$$\dot{x} = f(x) \quad (\text{B.3})$$

where $x \in \mathbf{R}^n$, $f : \mathbf{R}^n \rightarrow \mathbf{R}^n$.

Definition 1 [2] Let I be an interval of \mathbf{R} . A function $\phi : I \rightarrow \mathbf{R}^n$ is said to be a *Carathéodory solution* of (B.3) on I if $\phi(t)$ is absolutely continuous and $\frac{d\phi(t)}{dt} = f(\phi(t))$ for almost every $t \in I$.

Definition 3 [2] A function $V : \mathbf{R}^n \rightarrow \mathbf{R}$ is said to be *nonpathological* if it is locally Lipschitz continuous and for every absolutely continuous function $\phi : I \rightarrow \mathbf{R}^n$ and for almost every $t \in I$, the set $\partial_C V(\phi(t))$ is a subset of an affine subspace orthogonal to $\dot{\phi}(t)$, where $\partial_C V(x)$ denotes the Clarke gradient of real function V at point x .

Characterisation of Clarke gradient (Equation A.11 in [1]) If V is Lipschitz continuous, by Rademacher's theorem, its gradient $\nabla V(x)$ exists almost everywhere. Let N be the subset of \mathbf{R}^n where the gradient does not exist. It is possible to characterise Clarke generalised gradient as

$$\partial_C V(x) = \text{Co} \left\{ \lim_{x_i \rightarrow x} \nabla V(x_i), x_i \rightarrow x, x_i \notin N \cup \Omega \right\}$$

where Ω is any null measure set.

Definition 4 [2] Let $V : \mathbf{R}^n \rightarrow \mathbf{R}$ be a nonpathological function and Let (B.3) be given. Let

$$A_V = \{x \in \mathbf{R}^n : p_1 \cdot f(x) = p_2 \cdot f(x) \quad \forall p_1, p_2 \in \partial_C V(x)\}$$

if $x \in A_V$, the *nonpathological derivative of the map V with respect to (B.3) at x* is defined by

$$\dot{\bar{V}}_f(x) = p \cdot f(x)$$

where p is any vector in $\partial_C V(x)$.

Definition 5 [2] A vector field $f : \mathbf{R}^n \rightarrow \mathbf{R}^n$ is said to have the *solutions closure property* if for any sequence ϕ_n of solutions of (B.3) such that $\phi_n \rightarrow \phi$ uniformly on compact subsets of \mathbf{R} , one has that also ϕ is a solution of (B.3).

Definition 6 [2] A set M is said to be *weakly invariant* for (B.3) if for any $x_0 \in M$ there exists a $\phi \in S_{x_0}$, where S_{x_0} denotes the set of maximal solutions of (B.3) with initial condition $x(0) = x_0$, such that $\phi(t) \in M$ for all $t \geq 0$.

Proposition 3 [2] Assume that the vector field f has the solutions closure property. Let $V : \mathbf{R}^n \rightarrow \mathbf{R}$ be positive definite, nonpathological and radially unbounded. Let A_V be defined as in Definition 4 and assume

$$\dot{\bar{V}}_f(x) \leq 0, \quad \forall x \in A_V$$

Let $Z_f^V = \{x \in A_V : \dot{\bar{V}}_f(x) = 0\}$ and let M be the largest weakly invariant subset of Z_f^V . Then for any x_0 and any $\phi \in S_{x_0}$

$$\lim_{t \rightarrow +\infty} \text{dist}(\phi(t), M) = 0$$

B.3 Matrix Analysis

Theorem 4.2.2 (Rayleigh-Ritz) [18] Let $A \in M_n$ be Hermitian, and let the eigenvalues of A be ordered as

$$\lambda_{\min} = \lambda_1 \leq \lambda_2 \leq \cdots \leq \lambda_{n-1} \leq \lambda_n = \lambda_{\max}$$

Then

$$\lambda_1 x^* x \leq x^* A x \leq \lambda_n x^* x \quad \forall \quad x \in \mathbf{C}^n$$

$$\lambda_{\max} = \lambda_n = \max_{x \neq 0} \frac{x^* A x}{x^* x} = \max_{x^* x = 1} x^* A x$$

$$\lambda_{\min} = \lambda_1 = \min_{x \neq 0} \frac{x^* A x}{x^* x} = \min_{x^* x = 1} x^* A x$$

Theorem 6.1.1 (Geršgorin) [18] Let $A = [a_{ij}] \in M_n$, and Let

$$R'_i(A) = \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}|, \quad 1 \leq i \leq n$$

denote the *deleted absolute row sums* of A . Then all the eigenvalues of A are located in the union of n discs

$$\bigcup_{i=1}^n \{z \in \mathbf{C} : |z - a_{ii}| \leq R'_i(A)\} \equiv G(A)$$

Furthermore, if a union of k of these n discs forms a connected region that is disjoint from all the remaining $n - k$ discs, then there are precisely k eigenvalues of A in this region.

Definition 6.1.9 [18] Let $A = [a_{ij}] \in M_n$. The matrix A is said to be *diagonally dominant* if

$$|a_{ii}| \geq \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}| = R'_i, \quad i = 1, \dots, n$$

It is said to be *strictly diagonally dominant* if

$$|a_{ii}| > \sum_{\substack{j=1 \\ j \neq i}}^n |a_{ij}| = R'_i, \quad i = 1, \dots, n$$

Theorem 6.1.10 [18] Let $A = [a_{ij}] \in M_n$ be strictly diagonally dominant. Then

1. A is invertible.
2. If all main diagonal entries of A are positive, then all the eigenvalues of A have positive real parts.
3. If A is Hermitian and all main diagonal entries of A are positive, then all the eigenvalues of A are real and positive.

Bibliography

- [1] A. Bacciotti and F. Ceragioli. Nonsmooth optimal regulation and discontinuous stabilization. *Abstract and Applied Analysis*, 20:1159–1195, 2003.
- [2] A. Bacciotti and F. Ceragioli. Nonpathological Lyapunov functions and discontinuous carathéodory systems. *Automatica*, 42:453–458, 2006.
- [3] N. Bambos, S. C. Chen, and G. J. Pottie. Channel access algorithms with active link protection for wireless communication networks with power control. *IEEE/ACM Transactions on Networking*, 8(5):583–597, October 2000.
- [4] D. P. Bertsekas. *Nonlinear Programming*. Athena Scientific, Belmont, Massachusetts, 1999.
- [5] D. P. Bertsekas and R. Gallager. *Data Networks*. Prentice Hall, New Jersey, USA, 1992.
- [6] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press, Cambridge, UK, 2004.
- [7] L. Chen, S. H. Low, and J. C. Doyle. Joint congestion control and media access control design for ad hoc wireless networks. In *Proceedings of IEEE INFOCOM*, pages 2212–2222, March 2005.
- [8] M. Chiang, S. H. Low, A. R. Calderbank, and J. C. Doyle. Layering as optimization decomposition: A mathematical theory of network architectures. In *Proceedings of the IEEE*, volume 95, pages 255–312, January 2007.

- [9] M. Chiang, S. Zhang, and P. Hande. Distributed rate allocation for inelastic flows: Optimization frameworks, optimality conditions, and optimal algorithms. In *Proceedings of IEEE INFOCOM 2005*, volume 4, pages 2679–2690, March 2005.
- [10] R. L. Cruz. A calculus for network delay, part I: Network elements in isolation. *IEEE Transactions on Information Theory*, 37(1):114–131, January 1991.
- [11] R. L. Cruz. A calculus for network delay, part II: Network analysis. *IEEE Transactions on Information Theory*, 37(1):132–141, January 1991.
- [12] R. L. Cruz and A. V. Santhanam. Optimal routing, link scheduling and power control in multihop wireless networks. In *Proceedings of IEEE INFOCOM*, pages 702–711, April 2003.
- [13] T. ElBatt and A. Ephremides. Joint scheduling and power control for wireless ad hoc networks. *IEEE Transactions on Wireless Communications*, 3(1):74–85, January 2004.
- [14] A. V. Fiacco. *Introduction to Sensitivity and Stability Analysis in Nonlinear Programming*. Academic Press, New York, USA, 1983.
- [15] G. J. Foschini and Z. Miljanic. A simple distributed autonomous power control algorithm and its convergence. *IEEE Transactions on Vehicular Technology*, 42(4):641–646, November 1993.
- [16] A. Goldsmith. *Wireless Communications*. Cambridge University Press, USA, 2005.
- [17] P. Gupta and P. R. Kumar. The capacity of wireless networks. *IEEE Transactions on Information Theory*, 46(2):388–404, March 2000.
- [18] R. A. Horn and C. R. Johnson. *Matrix Analysis*. Cambridge University Press, USA, 1990.

- [19] S. Jahromizadeh and V. Rakocevic. Rate control for delay-sensitive traffic in multihop wireless networks. In *Proceedings of the 4th ACM Workshop on Performance Monitoring and Measurement of Heterogeneous Wireless and Wired Networks*, pages 99–106, Tenerife, Canary Islands, Spain, October 2009.
- [20] S. Jahromizadeh and V. Rakocevic. Joint rate control and scheduling for delay-sensitive traffic in multihop wireless networks. In *Proceedings of the IEEE 73rd Vehicular Technology Conference (VTC Spring)*, pages 1–5, Budapest, Hungary, May 2011.
- [21] H. K. Khalil. *Nonlinear Systems*. Prentice Hall, New Jersey, USA, 2002.
- [22] A. Lakshmikantha, C. L. Beck, and R. Srikant. Robustness of real and virtual queue-based active queue management schemes. *IEEE/ACM Transactions on Networking*, 13(1):81–93, Feb 2005.
- [23] J. Lee, R.R. Mazumdar, and N.B. Shroff. Non-convex optimization and rate control for multi-class services in the internet. *IEEE/ACM Transactions on Networking*, 13(4):827–840, 2005.
- [24] Y. Li, M. Chiang, and A. R. Calderbank. Congestion control in networks with delay sensitive traffic. In *Proceedings of IEEE GLOBECOM*, pages 2746–2751, 2007.
- [25] Z. Li and B. Li. Improving throughput in multihop wireless networks. *IEEE Transactions on Vehicular Technology*, 55(3):762–773, May 2006.
- [26] X. Lin and N. B. Shroff. Joint rate control and scheduling in multihop wireless networks. In *Proceedings of IEEE 43rd Conference on Decision and Control*, pages 1484–1489, Atlantis, Paradise Island, Bahamas, December 2004.
- [27] X. Lin and N. B. Shroff. The impact of imperfect scheduling on cross-layer congestion control in wireless networks. *IEEE/ACM Transactions on Networking*, 14(2):302–315, April 2006.

- [28] X. Lin and N. B. Shroff. Utility maximization for communication networks with multipath routing. *IEEE Transactions on Automatic Control*, 51(5):766–781, May 2006.
- [29] X. Lin, N. B. Shroff, and R. Srikant. A tutorial on cross-layer optimization in wireless networks. *IEEE Journal on Selected Areas in Communications*, 24(8):1452–1463, August 2006.
- [30] S. H. Low and D. E. Lapsley. Optimization flow control–I: Basic algorithm and convergence. *IEEE/ACM Transactions on Networking*, 7(6):861–874, December 1999.
- [31] I. Menache and A. Ozdaglar. *Network Games: Theory, Models, and Dynamics*. Synthesis Lectures on Communication Networks. Morgan and Claypool, USA, 2011.
- [32] F. Paganini, Z. Wang, J. C. Doyle, and S. H. Low. Congestion control for high performance, stability, and fairness in general networks. *IEEE/ACM Transactions on Networking*, 13(1):43–56, February 2005.
- [33] P. Santi. *Topology Control in Wireless Ad Hoc and Sensor Networks*. John Wiley and Sons, Chichester, England, 2005.
- [34] S. Shenker. Fundamental design issues for the future internet. *IEEE Journal on Selected Areas in Communications*, 13(7):1176–1188, September 1997.
- [35] R. Srikant. *The Mathematics of Internet Congestion Control*. Birkhäuser, Boston, USA, 2004.
- [36] S. Stidham. Pricing and congestion management in a network with heterogeneous users. *IEEE Transactions on Automatic Control*, 49(6):976–981, June 2004.
- [37] L. Trajkovic and S. J. Golestani. Congestion control for multimedia services. *IEEE Network*, 6(5):20–26, September 1992.

- [38] A. Tsirigos and Z. J. Haas. Multipath routing in the presence of frequent topological changes. *IEEE Communications Magazine*, 39(11):132–138, November 2001.
- [39] B. Wydrowski and M. Zukerman. QoS in best-effort networks. *IEEE Communications Magazine*, pages 44–49, December 2002.
- [40] T. Yoo, E. Setton, X. Zhu, A. Goldsmith, and B. Girod. Cross-layer design for video streaming over wireless ad hoc networks. In *Proceedings of IEEE 6th Workshop on Multimedia Signal Processing*, pages 99–102, Sienna, Italy, October 2004.
- [41] X. Zhu and B. Girod. Distributed rate allocation for multi-stream video transmission over ad hoc networks. In *Proceedings of IEEE International Conference on Image Processing*, pages II– 157–160, September 2005.